



## 저작자표시-비영리-변경금지 2.0 대한민국

이용자는 아래의 조건을 따르는 경우에 한하여 자유롭게

- 이 저작물을 복제, 배포, 전송, 전시, 공연 및 방송할 수 있습니다.

다음과 같은 조건을 따라야 합니다:



저작자표시. 귀하는 원저작자를 표시하여야 합니다.



비영리. 귀하는 이 저작물을 영리 목적으로 이용할 수 없습니다.



변경금지. 귀하는 이 저작물을 개작, 변형 또는 가공할 수 없습니다.

- 귀하는, 이 저작물의 재이용이나 배포의 경우, 이 저작물에 적용된 이용허락조건을 명확하게 나타내어야 합니다.
- 저작권자로부터 별도의 허가를 받으면 이러한 조건들은 적용되지 않습니다.

저작권법에 따른 이용자의 권리는 위의 내용에 의하여 영향을 받지 않습니다.

이것은 [이용허락규약\(Legal Code\)](#)을 이해하기 쉽게 요약한 것입니다.

[Disclaimer](#)

Doctoral Thesis

Learning based Wireless Communications with  
Energy Harvesting and Robot Vision Systems

Myeungun Kim

Department of Electrical Engineering

Graduate School of UNIST

2020

# Learning based Wireless Communications with Energy Harvesting and Robot Vision Systems

Myeungun Kim

Department of Electrical Engineering

Graduate School of UNIST

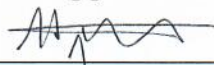
# Learning based wireless communications with energy harvesting and robot vision systems

A dissertation  
submitted to the Graduate School of UNIST  
in partial fulfillment of the  
requirements for the degree of  
Doctor of Philosophy

Myeungun Kim

6/12/2020

Approved by



---

Advisor

Hyun Jong Yang

# Learning based wireless communications with energy harvesting and robot vision systems

Myeungun Kim

This certifies that the dissertation of Myeungun Kim is approved.

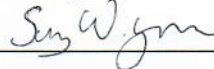
6/12/2020

signature



Advisor: Hyun Jong Yang

signature



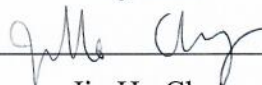
Sung Whan Yoon

signature



Youngbin Im

signature



Jin Ho Chung

signature



Young Chul Choi

## Abstract

From self-driving cars to smartphones essential to our lives, many types of the electronic devices and computers handle intelligently our work. Thanks to the ‘things’ that have become smarter, our lives have become more pleasant and faster, and literally easier. One of the big reasons we can live in such an environment is ‘machine learning’. It is a technology that allows a machine to acquire new knowledge by learning through a huge amount of data, just like a person learns. Machine learning is one of the most important topics in many industries and researches these days. It is no exaggeration to say that machine learning is used in almost every field. Its application to (1) wireless communications and (2) computer vision based robotics are also essential.

Learning based communication system has the following possibilities: (1) Unlike communication theory, real communication systems are non-linear. For this reason, deep learning-based communication systems may be more suitable for specific hardware configurations and channel optimization. (2) One of the great features of a communication system is that various signal processing functions (e.g., Coding, modulation, detection) are separated into several blocks. Rather than optimizing the performance of individual blocks, a machine learning-based end-to-end communication system can perform better. Because of these possibilities, machine learning is being applied to a wide range of communication systems such as heterogeneous access technology, cognitive radio, and resource allocation. In this dissertation, we propose a mathematical approach to the optimization problem of interference mitigation in a multi-cell network with and without energy harvesting. Also, we propose a recurrent neural network (RNN) based node selection algorithm for sensor networks with energy harvesting. Comparing the problem solving method of the former and the latter, the difference between the existing communication system and the learning based communication system can be clearly revealed.

Computer vision based robotics is a study that extracts meaningful information from an image or video and applies the information to a robot. In particular, as a result of applying machine learning to this field, various robots, such as autonomous vehicles, unmanned courier robots, and smart home robots, are being developed. The more studies on robots equipped with cameras, the more convenient our lives, but on the contrary, they can invade our privacy. That is, it is a double-edged sword. In this dissertation, we propose a method to protect our privacy while utilizing other visual information well (i.e., Simultaneous localization and mapping (SLAM)) by detecting faces in extreme low resolution images.



## Contents

I	Introduction . . . . .	1
1.1	Main Contributions . . . . .	1
II	Learning-based Wireless Communications with Energy Harvesting . . . . .	3
2.1	Downlink Beamforming in Small Cells with Scalar Information Exchange .	3
2.2	Min-SINR Maximization with DL SWIPT and UL WPCN in Multi-Antenna Interference Networks . . . . .	9
2.3	RNN-Based Node Selection for Sensor Networks with Energy Harvesting .	19
III	Learning-based Robot Vision Systems . . . . .	26
3.1	Privacy-Preserving Robot Vision with Anonymized Faces by Extreme Low Resolution . . . . .	26
	References . . . . .	39
	Acknowledgements . . . . .	44



## List of Figures

1	System model . . . . .	4
2	Achievable sum-rate/cell vs. SNR for $N_C = 3$ and $N_T = 2$ . . . . .	8
3	Proposed joint time switching SWIPT and WPCN protocol . . . . .	10
4	DL beamforming with energy harvesting and UL power allocation in multi-cell MISO networks . . . . .	11
5	Proposed max-min-SINR DL beamforming and UL PA design . . . . .	16
6	Minimum UL user rate vs. minimum DL user rate . . . . .	17
7	$\theta_{\min}, \theta_{\max}$ vs. number of iterations . . . . .	18
8	Superframe structure . . . . .	19
9	Basic structure of RNN . . . . .	20
10	System model . . . . .	21
11	UL packet format . . . . .	22
12	RNN structure of the proposed scheme . . . . .	22
13	Method to make labeled ground truth . . . . .	23
14	Training loss vs. Learning iterations . . . . .	24
15	Number of penalties vs. time . . . . .	24
16	Composition of the developed patrol robot system with privacy preserving face detection. . . . .	28

17	Dynamic resolution face detection architecture. . . . .	30
18	Our training data generation process for the example of the image with two faces. . . . .	32
19	Results comparing our proposed method with the results of the approach presented in Hu et al. for three different example images. . . . .	36
20	Face detection robot used in our experiments . . . . .	37
21	Comparison of the feature extraction results at various resolutions. . . . .	38

# I Introduction

From self-driving cars to smartphones essential to our lives, many types of the electronic devices and computers handle intelligently our work. Thanks to the ‘things’ that have become smarter, our lives have become more pleasant and faster, and literally easier. One of the big reasons we can live in such an environment is ‘machine learning’. It is a technology that allows a machine to acquire new knowledge by learning through a huge amount of data, just like a person learns. Machine learning is one of the most important topics in many industries and researches these days. It is no exaggeration to say that machine learning is used in almost every field. Its application to (1) wireless communications and (2) computer vision based robotics are also essential.

Learning based communication system has the following possibilities: (1) Unlike communication theory, real communication systems are non-linear. For this reason, deep learning-based communication systems may be more suitable for specific hardware configurations and channel optimization. (2) One of the great features of a communication system is that various signal processing functions (e.g., Coding, modulation, detection) are separated into several blocks. Rather than optimizing the performance of individual blocks, a machine learning-based end-to-end communication system can perform better. Because of these possibilities, machine learning is being applied to a wide range of communication systems such as heterogeneous access technology, cognitive radio, and resource allocation. In this dissertation, we propose a mathematical approach to the optimization problem of interference mitigation in a multi-cell network with and without energy harvesting. Also, we propose a recurrent neural network (RNN) based node selection algorithm for sensor networks with energy harvesting. Comparing the problem solving method of the former and the latter, the difference between the existing communication system and the learning based communication system can be clearly revealed.

Computer vision based robotics is a study that extracts meaningful information from an image or video and applies the information to a robot. In particular, as a result of applying machine learning to this field, various robots, such as autonomous vehicles, unmanned courier robots, and smart home robots, are being developed. The more studies on robots equipped with cameras, the more convenient our lives, but on the contrary, they can invade our privacy. That is, it is a double-edged sword. In this dissertation, we propose a method to protect our privacy while utilizing other visual information well (i.e., Simultaneous localization and mapping (SLAM)) by detecting faces in extreme low resolution images.

## 1.1 Main Contributions

This dissertation aims to provide the applications of machine learning to various fields, especially wireless communications and computer vision based robotics, where the main contributions are:

- Min-SINR Maximization with DL SWIPT and UL WPCN in Multi-Antenna Interference Networks.** Due to unpredictable future channel state, considering downlink (DL) and uplink (UL) jointly in multi-cell networks with energy harvesting is known to be very challenging. For resolving this issue, we have proposed a novel multi-cell communication and energy harvesting scheme, in which DL simultaneous wireless information and power transfer and UL wireless powered communication network concepts are jointly considered in [1]. Specifically, a DL beamforming with energy harvesting and UL power allocation (PA) scheme is proposed, where each cell is composed of a base station with multiple antennas and power-limited users each with single antenna. With the difficulty of optimizing mathematically DL and UL at the same time, it is considered to apply machine learning to a communication system as follows.
- RNN-Based Node Selection for Sensor Networks with Energy Harvesting.** Node selection problem in sensor networks is that a master node (MN) sequentially decides which slave node (SN) transmits UL data or receives DL data. However, the unpredictability of 1) future channel condition, 2) battery level of sensor devices, i.e., SN, and 3) packet deadlines of SNs makes the problem challenging. In [2], we propose an recurrent neural network (RNN) based node selection algorithm in pursuit of minimizing the transmission failures due to low battery level and exceeded UL/DL deadline. By processing sequential data with RNN, this algorithm takes into account battery level of sensor devices, UL/DL packet deadline as well as future channel information implicitly.
- Privacy-Preserving Robot Vision with Anonymized Faces by Extreme Low Resolution.** In the field of computer vision, privacy infringement is a serious social problem. In fact, the incident that the smart home camera memory is hacked occurs, and people's interest in protecting privacy has increased. In [3], we propose a learning based robot vision system which detects privacy-sensitive blocks, i.e., human face, from extreme low resolution (LR) images, and then dynamically enhances the resolution of only privacy-insensitive blocks, e.g., backgrounds. Keeping all the face blocks to be extreme LR of 15x15 pixels, we can guarantee that human faces are never at high resolution (HR) in any of processing or memory, thus yielding strong privacy protection even from cracking or backdoors.

## II Learning-based Wireless Communications with Energy Harvesting

### 2.1 Downlink Beamforming in Small Cells with Scalar Information Exchange

#### Introduction

The recent evolution of mobile wireless communications is being driven by the concepts of small cells [4], [5] and multiinput and multi-output (MIMO) [6–8]. The idea of small cells is to increase the infrastructure density by employing physically and functionally ‘small’ BSs within a macro cell coverage. With the cell densification, path loss from BSs to user equipment (UE) is expected to be significantly reduced, yielding higher data rate over the entire coverage area [8]. In addition, in MIMO systems, the spectral efficiency can be scaled by transmitting and receiving independent information on different spatial streams [6], [8], or the link reliability can be improved by employing beamforming techniques [7], [8]. It is, needless to say, that the spectral efficiency can be further enhanced if both the small cell and MIMO techniques are used concurrently.

Four small cell scenarios are discussed by the 3GPP standard body [4], [5] considering several different aspects: same and separate frequency for macro and small cells, indoor and outdoor deployment, and with and without a macro coverage. One of the major problems in small cells is intercell interference, which becomes more significant as the cell density increases.

In this paper, we address the intercell interference issue between small cells without consideration of the macro coverage; that is, the scenario with separate frequency deployment or the scenario without a macro coverage is considered. In particular, the scenario where the BSs are equipped with multiple antennas is addressed. We first revisit the two extreme beamforming strategies: 1) max-signal-to-noise-ratio (SNR) and 2) min-generating-interference (GI) schemes. In the max-SNR scheme, the maximum transmit combining (MTC) beamforming scheme is considered only to maximize the desired channel gain at each cell. On the other hand, in the min-GI scheme, each BS designs its beamforming vector such that the interference it generates is minimized as done in [9], [10]. These two schemes require only local channel state information (CSI) at transmitter (CSIT) assuming the time-division duplexing (TDD) channel reciprocity, i.e., each BS has the knowledge of the incoming and outgoing channels and is ignorant of the channels between the BSs and UE in other cells.

Then, we propose a beamforming scheme that minimizes the weighted-GI (WGI) to in pursuit of further enhancing the achievable rate. The weight coefficients are determined to take into account both the desired channel gain and generating interference, thereby balancing between the two extreme philosophies – egoism and altruism.

Specifically, the beamforming vector at each BS is designed such that it minimizes the WGI, where the weight coefficients are determined according to the signal-to-interference-and-noise-

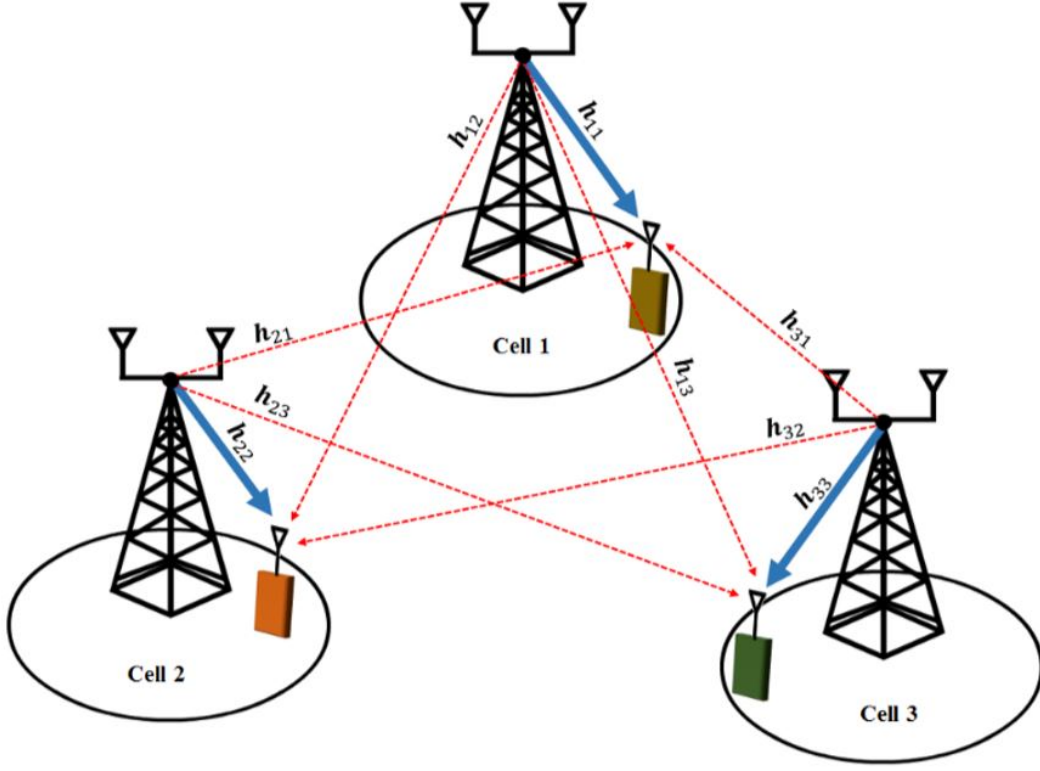


Figure 1: System model

ratio (SINR) at the other cells. Thus, the proposed scheme requires an additional information exchange between the BSs for a few scalar values such as the SINRs. The scheme is still feasible, since in all the four 3GPP scenarios, a certain amount of information is allowed to be exchanged between small cell BSs via X-2 interface [4], [5]. Simulation results show that there exists a SINR cross over a point from which the performance of the min-GI scheme surpasses the performance of the max-SNR scheme. In addition, it is shown that the proposed min-WGI scheme always outperforms the other schemes in all SINR regime.

The remainder of this paper is organized as follows. Section II describes the system model, and Section III presents the three different methods to design the beamforming vectors. Section IV provides numerical simulation results, and Section V concludes the paper.

## System model

As shown in Fig. 1, the small cell scenario without a macro cell coverage is considered. For simplicity of notation, only a single UE is depicted. It is straightforward to extend the study to the multi-user case with the use of orthogonal frequency division multiplexing [11].

It is assumed that each small cell BS has  $N_T$  antennas, while each UE has a single antenna. The number of small cells is denoted by  $N_C$ . The channel vector from the  $i$ -th BS to the UE in the  $j$ -th cell is denoted by  $\mathbf{h}_{ij} \in \mathbb{C}^{1 \times N_T}$ ; that is, the first letter of the subscript denotes the source and the second denotes the destination. Block fading is assumed. TDD and the channel

reciprocity is assumed. It is also assumed that each BS has the knowledge of the channels from itself, i.e., the  $i$ -th BS has the information of  $\mathbf{h}_{ij}$ ,  $j = 1, \dots, N_C$ .

The beamforming vector at the  $i$ -th BS is denoted by  $\mathbf{w}_i \in \mathbb{C}^{N_T \times 1}$ , where  $\|\mathbf{w}_i\|^2 = 1$ . The received signal at the UE in the cell is written by

$$y_i = \underbrace{\mathbf{h}_{ii}\mathbf{w}_i x_i}_{\text{desired signal}} + \underbrace{\sum_{k=1, k \neq i}^{N_C} \mathbf{h}_{ki}\mathbf{w}_k x_k}_{\text{intercell interference}} + z_i, \quad (1)$$

where  $x_i$  is the unit-variance transmit symbol, and  $z_i$  is the additive white Gaussian noise (AWGN) at the UE in the  $i$ -th cell. Thus, the corresponding SINR is expressed by

$$\text{SINR}_i = \frac{|\mathbf{h}_{ii}\mathbf{w}_i|^2}{\sum_{k=1, k \neq i}^{N_C} |\mathbf{h}_{ki}\mathbf{w}_k|^2 + N_0}, \quad (2)$$

and the total achievable sum-rate is given by

$$R = \sum_{i=1}^{N_C} \log(1 + \text{SINR}_i). \quad (3)$$

## Beamforming design

- Max-SNR

In the max-SNR scheme, the beamforming vector is designed such that the channel gain of the desired channel, i.e.,  $\|\mathbf{h}_{ii}\mathbf{w}_i\|^2$ , is maximized such that

$$\mathbf{w}_i^{\text{max-SNR}} = \max_{\mathbf{w}} |\mathbf{h}_{ii}\mathbf{w}|^2, \text{ s.t. } \|\mathbf{w}_i\|^2 = 1. \quad (4)$$

Thus, the solution for the max-SNR is a simple maximum ratio transmit (MRT) scheme [12] given by

$$\mathbf{w}_i^{\text{max-SNR}} = \frac{\mathbf{h}_{ii}^H}{\|\mathbf{h}_{ii}\|}. \quad (5)$$

Note that this solution is designed independently of the interference channels, and hence, it does not change the strength of the interference channels in an average sense.

- Min-Generating-Interference

With the local CSIT, each BS can calculate the amount of interference that it generates from:

$$\Delta_i = \sum_{j=1, j \neq i}^{N_C} |\mathbf{h}_{ij}\mathbf{w}_i|^2. \quad (6)$$

Therefore, to minimize the GI, the beamforming vector is designed such that

$$\mathbf{w}_i^{\text{min-interf}} = \min_{\mathbf{w}} \Delta_i(\mathbf{w}), \text{ s.t. } \|\mathbf{w}_i\|^2 = 1. \quad (7)$$

Note that  $\Delta_i$  can be expressed as

$$\Delta_i = \left\| \begin{bmatrix} \mathbf{h}_{i1} \\ \vdots \\ \mathbf{h}_{i(i-1)} \\ \mathbf{h}_{i(i+1)} \\ \vdots \\ \mathbf{h}_{iN_C} \end{bmatrix} \mathbf{w}_i \right\|^2 = \|\mathbf{G}_i \mathbf{w}_i\|^2. \quad (8)$$

If we denote the singular value decomposition of the matrix  $\mathbf{G}_i \in \mathbb{C}^{(N_C-1) \times N_T}$  as

$$\mathbf{G}_i = \mathbf{U}_i \mathbf{\Sigma}_i \mathbf{V}_i^H, \quad (9)$$

where  $\mathbf{U}_i \in \mathbb{C}^{(N_C-1) \times \bar{N}}$  and  $\mathbf{V}_i \in \mathbb{C}^{\bar{N} \times N_T}$  consist of orthogonal columns, and  $\mathbf{\Sigma}_i \in \mathbb{R}^{\bar{N} \times \bar{N}}$  is a diagonal matrix composed of singular values of  $\mathbf{G}_i$ . Here  $\bar{N} = \min(N_C - 1, N_T)$ . Now, the solution for the problem (7) is given by

$$\mathbf{w}_i^{\min\text{-interf}} = \mathbf{v}_{i,\bar{N}}, \quad (10)$$

where  $\mathbf{v}_{i,\bar{N}}$  is the  $\bar{N}$ -th column of  $\mathbf{V}_i$ , which is associated with the minimum singular value of the interference matrix  $\mathbf{G}_i$ .

*Remark 1.* Since the beamforming vector is designed only to minimize the GI, the desired channel gains do not change in an average sense.

*Remark 2.* In [9], [10], it was shown that the min GI scheme suffices to achieve the optimal degrees-of-freedom (DoF) in multicell networks. However, it was also shown that the minGI scheme is not optimal in terms of the achievable sum-rate. Furthermore, it is an open problem to find an optimal beamforming scheme achieving the maximum achievable sum-rate with local CSIT and/or even with limited information exchange between the BSs.

- Min Weighted Generating-Interference

For the structured beamforming scheme, it is assumed that scalar values can be exchanged between the small cell BSs. We first define the weighted generating-interference (WGI) as follows:

$$\Omega_i = \sum_{j=1, j \neq i}^{N_C} \beta_{ij} |\mathbf{h}_{ij} \mathbf{w}_i|^2, \quad (11)$$

where  $\beta_{ij} \geq 0$  and

$$\sum_{j=1, j \neq i}^{N_C} \beta_{ij} = 1. \quad (12)$$

Here, the interference weight  $\beta_{ij}$  accounts for the relative emphasis on each interference channel. From the fact that the sum-rate of MIMO systems with multiple spatial streams



is maximized if more emphasis is put on the spatial streams with higher SNRs, i.e., water-filling power allocation [12], we propose to determine the weights  $\beta_{ij}$  such that it is proportional to the corresponding SINRs. Formally,  $\beta_{ij}$  is determined as

$$\beta_{ij} = \frac{\rho_{ij}}{\sum_{k=1, k \neq i}^{N_C} \rho_{ik}}, j = 1, \dots, i-1, i+1, \dots, N_C, \quad (13)$$

where

$$\rho_{ij} = \frac{\|\mathbf{h}_{jj}\|^2}{\sum_{k=1, k \neq j}^{N_C} \|\mathbf{h}_{kj}\|^2 + N_0}. \quad (14)$$

Following the same footsteps of the min-GI scheme, we find the beamforming vector that minimizes the WGI as

$$\mathbf{w}_i^{\text{min-WGI}} = \mathbf{v}_{i, \bar{N}}, \quad (15)$$

where  $\mathbf{v}_{i, \bar{N}}$  is the right singular vector associated with the minimum singular value of the weighted interference matrix

$$\mathbf{G}_i = \begin{bmatrix} \sqrt{\beta_{i1}} \mathbf{h}_{i1} \\ \vdots \\ \sqrt{\beta_{i(i-1)}} \mathbf{h}_{i(i-1)} \\ \sqrt{\beta_{i(i+1)}} \mathbf{h}_{i(i+1)} \\ \vdots \\ \sqrt{\beta_{iN_C}} \mathbf{h}_{iN_C} \end{bmatrix} \in \mathbb{C}^{(N_C-1) \times N_T}. \quad (16)$$

*Remark 3.* [Decoupling of the beamforming design and SNR calculation]. In fact, the actual SINRs depend on the design of beamforming vectors as in (2). Hence, the determination of  $\beta_{ij}$  is coupled with the SINR calculation, and thus the beamforming design cannot be done independently of other cells. We propose to use the pre-processing SINRs  $\rho_{ij}$  given by (14) to decouple the beamforming design procedures at all the cells. Simulation results shall show that even with this suboptimality, the achievable sum-rate is significantly improved compared to the min-GI scheme.

*Remark 4.* [Scalar information exchange between small cell BSs]. Note that to design  $\beta_{ij}$ ,  $j = 1, \dots, i-1, i+1, \dots, N_C$ , the  $i$ -th BS needs to have the knowledge about the SINRs in the other cells,  $\rho_{ij}$ ,  $j = 1, \dots, i-1, i+1, \dots, N_C$ . The relevant channels  $\mathbf{h}_{jj}$ ,  $j = 1, \dots, i-1, i+1, \dots, N_C$ , are not subject to local CSIT at  $i$ -th BS's viewpoint, and hence, these scalar values need to be exchanged between the BSs. This exchange can be easily done via a low-rate direct interface between the BSs, such as the X2 interface in the 3GPP small cell systems [4], [5].

## Numerical results

The achievable sum-rates of the three beamforming schemes are compared under Rayleigh fading environment. For comparison, the baseline ‘Random’ scheme is also considered, in which all the beamforming vectors are randomly chosen.

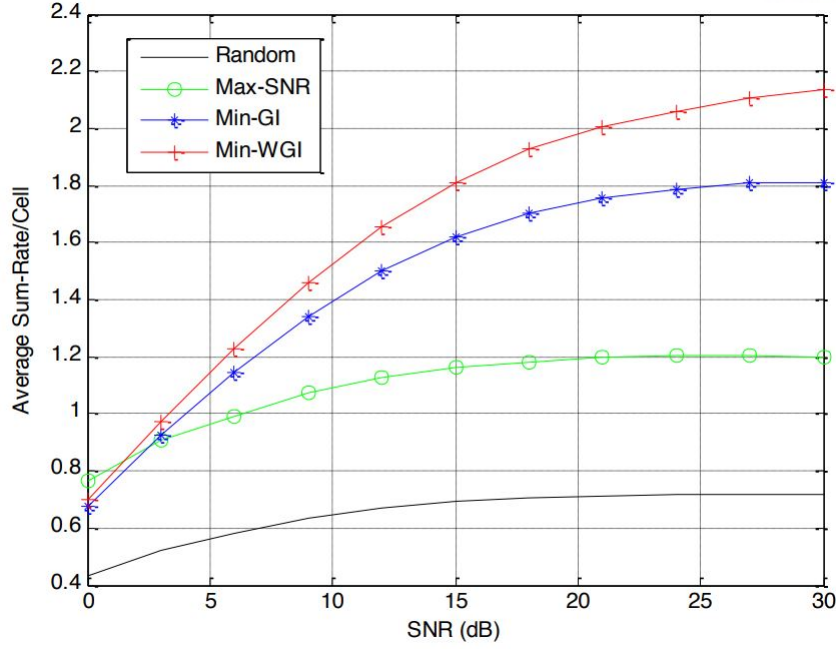


Figure 2: Achievable sum-rate/cell vs. SNR for  $N_C = 3$  and  $N_T = 2$

Fig. 2 shows the average achievable sum-rate per cell versus SNR for  $N_C = 3$  and  $N_T = 2$ . For the noise-limited regime, i.e., low SNR regime, the max-SNR scheme shows relatively high performance than the others, since the interference is predominated by the noise. On the other hand, as the SNR increases, the efforts to reduce the GI provide sum-rate gain over the max-SNR scheme as seen from Fig. 2, since the sum-rate is dominated by the interference.

It is also seen that the min-WGI scheme exhibits notable sum-rate gain compared to the min-GI, as the min-WGI scheme also takes into account the desired channel gains as well as the GI. Note again that this sum-rate gain is achieved at the cost of additional low-rate scalar information exchange between the BSs.

## Conclusion

We have revisited two extreme beamforming schemes with local CSIT, in which the beamforming vector is designed based on the egoism and altruism philosophies, respectively. In addition, we have proposed a new beamforming scheme that minimizes the weighted generating-interference, where the weight coefficients are determined proportionally to the SINRs at the neighboring cells. The determination of these weight coefficients also takes into account the desired channel gains as well as the generating-interference, and hence, the scheme requires addition scalar information exchange between the small cell BSs. Through simulation results, we have shown that the proposed min-generating-interference scheme outperforms the other schemes for mid to high SNR regime.

## 2.2 Min-SINR Maximization with DL SWIPT and UL WPCN in Multi-Antenna Interference Networks

### Introduction

In dense cellular networks, one of the most critical problems that degrade the cell throughput is inter-cell interference (ICI) [13]. To mitigate the ICI in the downlink (DL) scenario, multi-antenna DL beamforming has been introduced. Assuming global channel state information (CSI), the DL beamforming vectors were optimized in the sense of maximizing the sumrate [14].

Recently, energy harvesting (EH) from ambient RF signals has received intense research interest to save energy consumption on battery-limited devices such as low power sensors and mobile devices. Several studies have shown the feasibility of EH systems based on ambient RF signals in practical environments [15], [16] and cellular network [17]. To utilize interference signals for energy harvesting, the concept of simultaneous wireless information and power transfer (SWIPT) has been extensively studied both theoretically and practically [18]. It becomes more crucial to harvest energy from ICI signals in wireless sensor networks or ultra dense networks, where the ICI is even comparable in strength to the desired signal. In DL cellular network based on global CSI [19] or local CSI [20], [21], multi-antenna beamforming or precoding at the base stations (BSs) can guarantee the data rate of information transfer as well as the amount of energy harvested at the users. In particular, the benefit of the SWIPT is emphasized in highly dense networks such as small cells, since strong ICI in such a network can be utilized as a source of EH [22]. However, these previous schemes merely focus on DL signal-to-interference-plus-noise ratios (SINRs) and the amount of energy harvested without any explicit consideration of the uplink (UL) SINRs.

Another framework ‘wireless powered communication network (WPCN)’ [23] has been proposed to consider UL information transmission, where each user is powered by the energy that it harvests from the DL signals in DL time slots. The authors of [23] optimized the time allocation for the DL wireless energy transfer and UL information transmission to maximize UL sum-throughput. This work has been extended to multiuser multi-input single-output (MU-MISO) [24] and MU multi-input multi-output (MIMO) [25] channels. Nevertheless, these WPCN schemes only consider energy transfer without information transmission in DL time slots, and thus they are not applicable to the case with the presence of DL data, which is highly likely in practical systems. In addition, these works considered only the single cell case, and the extension to the multicell scenario is non-trivial in managing the ICI both in the DL and UL time slots.

In this letter, we tackle a general communication and energy transfer scenario for multicell networks composed of BSs with multiple antennas and users each with a single antenna, where there exist both DL and UL data to exchange. In the proposed scheme, the concepts of SWIPT and WPCN are jointly considered for the DL and UL time slots, and the UL power allocation (PA) is also considered in the UL time slots. Specifically, in a DL time slot, while a user receives DL data, another user scavenges all the ambient signals transmitted by the BSs. In

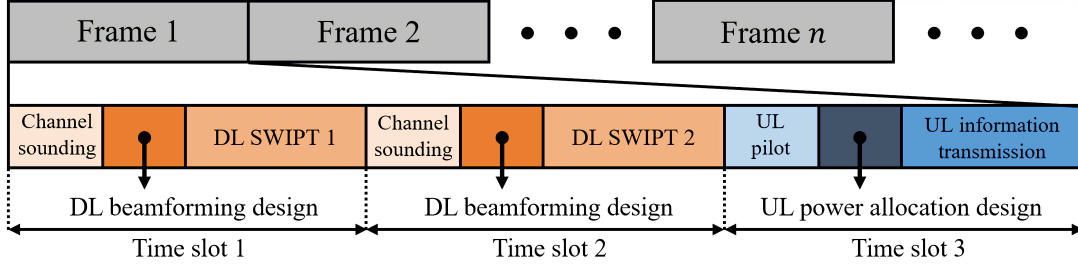


Figure 3: Proposed joint time switching SWIPT and WPCN protocol

the subsequent DL time slot, the role of the two users is swapped. For each DL time slot, beamforming vectors at the BSs are optimized. In the following UL time slot, the two users transmit UL data simultaneously with a proper UL PA by each BS. In fact, the coupled DL and UL problems are known to be challenging to solve since the optimality of the DL and UL parameter design at each time slot does not exist in time-varying channels. Therefore, we propose an efficient cascaded DL and UL design scheme maximizing the minimum DL and UL SINRs in pursuit of maximizing the rate fairness. Simulation results show that the proposed scheme not only achieves larger minimum DL and UL SINRs region but also exhibits much improved energy efficiency.

### System Model and Proposed Protocol

An  $N_C$ -cell network, each cell of which is composed of a BS with  $N_T$  antennas and two users each with a single antenna, is considered. The extension to the case with more users is trivial assuming each pair of users is orthogonalized with the other pairs of users based on multi-carrier modulation such as orthogonal frequency division multiplexing. We consider frame-based with a three-time slot protocol shown in Fig. 3. In the first time slot, user 1 receives the DL data while the other harvests energy from the DL signals transmitted by the BSs. In the second time slot, the two users switch their roles. In the third time slot, all the users transmit UL signals simultaneously using the energy harvested from the DL time slots. That is, ‘time switching’ SWIPT is assumed, which requires only a circuit switch for implementation. Compared to ‘power splitting’ SWIPT, where each user can split the received energy into two parts for EH and information decoding by an additional power splitting circuit at each RF chain, the time switching SWIPT requires relatively low implementation cost [18].

It is also assumed that each user is equipped with a battery such that the energy unused can be stored. The channel vector between the BS in cell  $i$  and user  $m$  in cell  $j$  at the  $n$ -th slot is denoted by  $\mathbf{h}_{i,(j,m)}^{[n]} \in \mathbb{C}^{N_T \times 1}$ ,  $i, j = 1, \dots, N_C$ ,  $m = 1, 2$ , and  $n = 1, 2, 3$ . It is assumed that the channel coefficients remain constant for a time slot and then change to another values randomly at the next time slot, i.e., quasi-static fading. Each BS is assumed to be able to acquire its incoming and outgoing channels through channel sounding and UL pilot signals at each time slot and share them with the other BSs via backhaul.

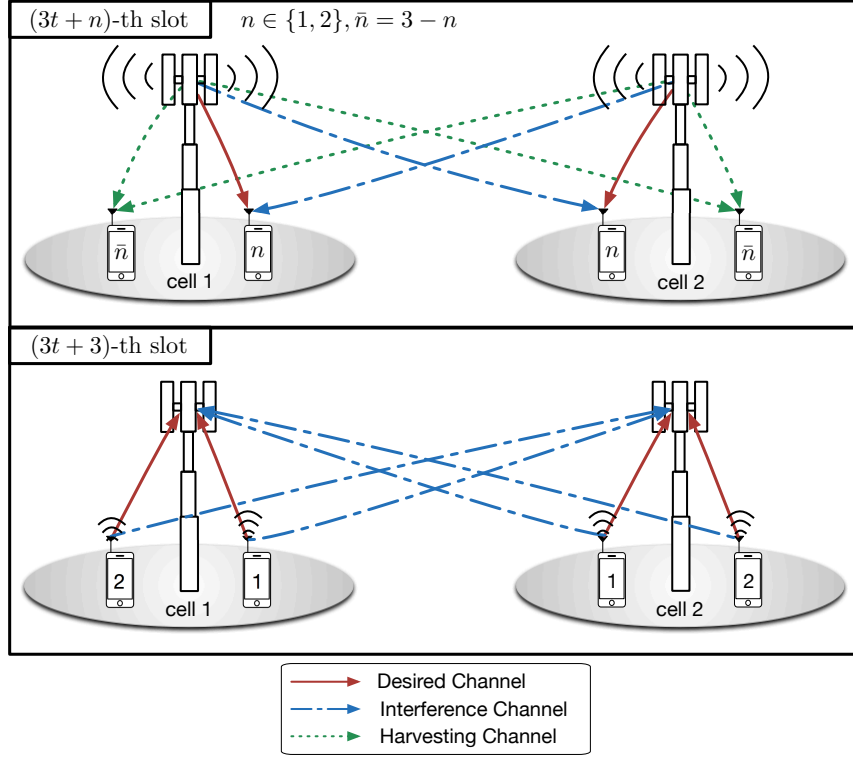


Figure 4: DL beamforming with energy harvesting and UL power allocation in multi-cell MISO networks

Fig. 4 shows the system model for  $N_C = 2$  and  $N_T = 3$ . At the  $n$ -th time slot,  $n = 1, 2$ , i.e., the DL time slot, the beamforming vector at the BS in cell  $j$  is denoted by  $\mathbf{w}_j^{[n]} \in \mathbb{C}^{N_T \times 1}$ , where  $\|\mathbf{w}_j^{[n]}\|^2 = 1$ . The DL SINR for user  $n$  in cell  $j$  is thus

$$\rho_{j,n}^{[n]} = \frac{P_j \left| \left( \mathbf{h}_{j,(j,n)}^{[n]} \right)^H \mathbf{w}_j^{[n]} \right|^2}{\sum_{k=1, k \neq j}^{N_C} P_k \left| \left( \mathbf{h}_{k,(j,n)}^{[n]} \right)^H \mathbf{w}_k^{[n]} \right|^2 + N_0}, \quad (17)$$

where  $P_j$  is the transmit power of the BS in cell  $j$ , and  $N_0$  is the variance of additive white Gaussian noise. In addition,  $\nu_{j,\bar{n}}$  denotes the energy harvested at user  $(3-n)$  in cell  $j$ , defined by

$$\nu_{j,\bar{n}} = \gamma_{j,\bar{n}} \sum_{k=1}^{N_C} \left| \left( \mathbf{h}_{k,(j,\bar{n})}^{[n]} \right)^H \mathbf{w}_k^{[n]} \right|^2, \quad (18)$$

where  $\bar{n} = 3 - n$ , and  $0 < \gamma_{j,\bar{n}} \leq 1$  denotes the EH efficiency, the ratio of the energy stored in the battery to the total received power of user  $\bar{n}$  in cell  $j$  [26], [18].

At the third time slot, i.e., the UL time slot, we assume a linear receiver with successive interference cancellation [27] based on the decoding order  $\boldsymbol{\pi}_j$  at the BS in cell  $j$ , where  $\boldsymbol{\pi}_j = [\pi_{(j,1)}, \pi_{(j,2)}]$ ,  $\pi_{(j,1)}, \pi_{(j,2)} \in \{1, 2\}$ . The user with index  $\pi_{(j,1)}$  is decoded first and the user with index  $\pi_{(j,2)}$  is decoded second after the subtraction of the interference due to user  $\pi_{(j,1)}$ 's signal.

The interference-plus-noise spatial covariance matrix of user  $\pi_{(j,m)}$  in cell  $j$  is given by

$$\mathbf{Z}_{j,\pi_{(j,1)}} = \underbrace{p_{j,\pi_{(j,2)}} \mathbf{h}_{j,(j,\pi_{(j,2)})}^{[3]} \left( \mathbf{h}_{j,(j,\pi_{(j,2)})}^{[3]} \right)^H}_{\text{intra-cell interference}} + \mathbf{C}_j + N_0 \mathbf{I}, \quad (19)$$

$$\mathbf{Z}_{j,\pi_{(j,2)}} = \mathbf{C}_j + N_0 \mathbf{I}, \quad (20)$$

where  $p_{j,\pi_{(j,m)}}$  is the UL transmit power from user  $\pi_{(j,m)}$  and

$$\mathbf{C}_j = \underbrace{\sum_{k=1, k \neq j}^{N_C} \sum_{m=1}^2 p_{k,\pi_{(k,m)}} \mathbf{h}_{j,(k,\pi_{(k,m)})}^{[3]} \left( \mathbf{h}_{j,(k,\pi_{(k,m)})}^{[3]} \right)^H}_{\text{inter-cell interference}}. \quad (21)$$

Denoting the receiver beamforming vector for user  $\pi_{(j,m)}$  by  $\mathbf{u}_{j,\pi_{(j,m)}} \in \mathbb{C}^{N_T \times 1}$ , where  $\|\mathbf{u}_{j,\pi_{(j,m)}}\|^2 = 1$ , we obtain the UL SINR as

$$\rho_{j,\pi_{(j,m)}}^{[3]} = \frac{p_{j,\pi_{(j,m)}} \left| \left( \mathbf{u}_{j,\pi_{(j,m)}} \right)^H \mathbf{h}_{j,(j,\pi_{(j,m)})}^{[3]} \right|^2}{\left( \mathbf{u}_{j,\pi_{(j,m)}} \right)^H \mathbf{Z}_{j,\pi_{(j,m)}} \mathbf{u}_{j,\pi_{(j,m)}}}. \quad (22)$$

The goal is to maximize the minimum DL and UL SINRs by optimizing  $\mathbf{w}_j^{[n]}$  and  $p_{j,n}$ ,  $j = 1, \dots, N_C$ ,  $n = 1, 2$ , which requires a joint optimization for the three time slots. To decouple the optimization at each time slot, we propose to design  $\mathbf{w}_j^{[n]}$  to maximize the minimum DL SINR at the first and second time slots while guaranteeing the minimum amount of harvested energy for UL information transmission at each user. In the third time slot, the UL power  $p_{j,n}$  is optimized to maximize the minimum UL SINR.

### Optimization of the DL beamforming

We begin with defining the augmented variable vector by  $\tilde{\mathbf{w}}^{[n]} \triangleq \left[ \left( \mathbf{w}_1^{[n]} \right)^T \dots \left( \mathbf{w}_{N_C}^{[n]} \right)^T \right]^T \in \mathbb{C}^{(N_C \cdot N_T) \times 1}$ . Then, each beamforming vector can be expressed by  $\mathbf{w}_j^{[n]} = \mathbf{E}_j \tilde{\mathbf{w}}^{[n]}$ , where  $\mathbf{E}_j \in \mathbb{C}^{N_T \times (N_C \cdot N_T)}$  consists of all zeros except for the  $((j-1)N_T + 1)$ -th to  $(jN_T)$ -th columns equal to the  $N_T \times N_T$  identity matrix. We can rewrite (17) as

$$\rho_{j,n}^{[n]} = \frac{(\tilde{\mathbf{w}}^{[n]})^H \mathbf{A}_{j,n}^{[n]} \tilde{\mathbf{w}}^{[n]}}{(\tilde{\mathbf{w}}^{[n]})^H \mathbf{B}_{j,n}^{[n]} \tilde{\mathbf{w}}^{[n]}}, \quad (23)$$

where

$$\mathbf{A}_{j,n}^{[n]} = P_j \mathbf{E}_j^H \mathbf{h}_{j,(j,n)}^{[n]} \left( \mathbf{h}_{j,(j,n)}^{[n]} \right)^H \mathbf{E}_j, \quad (24)$$

$$\mathbf{B}_{j,n}^{[n]} = \sum_{k=1, k \neq j}^{N_C} \left( P_k \mathbf{E}_k^H \mathbf{h}_{k,(j,n)}^{[n]} \left( \mathbf{h}_{k,(j,n)}^{[n]} \right)^H \mathbf{E}_k \right) + N_0 \mathbf{I}, \quad (25)$$

and  $\mathbf{I}$  is the  $(N_C \cdot N_T) \times (N_C \cdot N_T)$  identity matrix. We also rewrite (18) as

$$\nu_{j,\bar{n}} = \gamma_{j,\bar{n}} \sum_{k=1}^{N_C} \left| \left( \mathbf{h}_{k,(j,\bar{n})}^{[n]} \right)^H \mathbf{E}_k \tilde{\mathbf{w}}^{[n]} \right|^2. \quad (26)$$

The optimization problem at the  $n$ -th time slot,  $n = 1, 2$ , is formulated by

$$\mathcal{P}_{D1}^{[n]} : \max_{\tilde{\mathbf{w}}^{[n]}} \min_j \left\{ \frac{(\tilde{\mathbf{w}}^{[n]})^H \mathbf{A}_{j,n}^{[n]} \tilde{\mathbf{w}}^{[n]}}{(\tilde{\mathbf{w}}^{[n]})^H \mathbf{B}_{j,n}^{[n]} \tilde{\mathbf{w}}^{[n]}} \right\} \quad (27)$$

$$\text{s.t. } \nu_{j,\bar{n}} \geq \Lambda_{j,\bar{n}} - \Delta_{j,\bar{n}}, \quad (28)$$

$$\|\mathbf{E}_j \tilde{\mathbf{w}}^{[n]}\|^2 \leq 1, \quad j = 1, \dots, N_C. \quad (29)$$

In addition,  $\Lambda_{j,\bar{n}}$  is the minimum energy to be stored at user  $\bar{n}$  in cell  $j$  for UL transmit in the UL time slot, and  $\Delta_{j,\bar{n}}$  is the power left in the battery of user  $\bar{n}$  in cell  $j$ . By controlling  $\Lambda_{j,\bar{n}}$ , the proposed protocol can cover general scenarios with different EH constraints, while minimizing unnecessary energy consumption at the BSs. Due to the constraint (28), the feasible set of the problem  $\mathcal{P}_{D1}^{[n]}$  becomes non-convex, which in general is difficult to solve. To convert the problem into a convex form, the maximization of the minimum of DL SINRs in (27) is rewritten employing one additional variable  $\theta$  by

$$\mathcal{P}_{D2}^{[n]} : \max_{\tilde{\mathbf{w}}^{[n]}, \theta} \theta \quad (30)$$

$$\text{s.t. } (\tilde{\mathbf{w}}^{[n]})^H \mathbf{A}_{j,n}^{[n]} \tilde{\mathbf{w}}^{[n]} \geq \theta (\tilde{\mathbf{w}}^{[n]})^H \mathbf{B}_{j,n}^{[n]} \tilde{\mathbf{w}}^{[n]}, \quad (31)$$

$$(\tilde{\mathbf{w}}^{[n]})^H \mathbf{C}_{j,\bar{n}}^{[n]} \tilde{\mathbf{w}}^{[n]} \geq \frac{\Lambda_{j,\bar{n}} - \Delta_{j,\bar{n}}}{\gamma_{j,\bar{n}}}, \quad (32)$$

$$(\tilde{\mathbf{w}}^{[n]})^H \mathbf{D}_j \tilde{\mathbf{w}}^{[n]} \leq 1, \quad j = 1, \dots, N_C, \quad (33)$$

where

$$\mathbf{C}_{j,\bar{n}}^{[n]} = \sum_{k=1}^{N_C} \mathbf{E}_k^H \mathbf{h}_{k,(j,\bar{n})}^{[n]} \left( \mathbf{h}_{k,(j,\bar{n})}^{[n]} \right)^H \mathbf{E}_k, \quad (34)$$

$$\mathbf{D}_j = \mathbf{E}_j^H \mathbf{E}_j. \quad (35)$$

The problem  $\mathcal{P}_{D2}^{[n]}$  is still in a non-convex form due to the non-convex constraints (31) and (32). Therefore, defining a rank-1 variable matrix by  $\mathbf{W}^{[n]} \triangleq \tilde{\mathbf{w}}^{[n]} (\tilde{\mathbf{w}}^{[n]})^H \in \mathbb{C}^{(N_C \cdot N_T) \times (N_C \cdot N_T)}$ , for given cost value  $\theta$ , we consider the following feasibility problem:

$$\mathcal{P}_{D3}^{[n]}(\theta) : \text{Find } \mathbf{W}^{[n]} \quad (36)$$

$$\text{s.t. } \text{tr} \left( \mathbf{A}_{j,n}^{[n]} \mathbf{W}^{[n]} \right) \geq \theta \text{tr} \left( \mathbf{B}_{j,n}^{[n]} \mathbf{W}^{[n]} \right), \quad (37)$$

$$\text{tr} \left( \mathbf{C}_{j,\bar{n}}^{[n]} \mathbf{W}^{[n]} \right) \geq \frac{\Lambda_{j,\bar{n}} - \Delta_{j,\bar{n}}}{\gamma_{j,\bar{n}}}, \quad (38)$$

$$\text{tr} \left( \mathbf{D}_j \mathbf{W}^{[n]} \right) \leq 1, \quad j = 1, \dots, N_C, \quad (39)$$

$$\mathbf{W}^{[n]} \succeq \mathbf{0}, \quad (40)$$

$$\text{rank} \left( \mathbf{W}^{[n]} \right) = 1, \quad (41)$$

where  $\text{tr}(\cdot)$  denotes the trace operation and  $\mathbf{W}^{[n]} \succeq \mathbf{0}$  implies that  $\mathbf{W}^{[n]}$  is positive semidefinite.

The feasibility problem  $\mathcal{P}_{\text{D3}}^{[n]}(\theta)$  without the constraint (41) is a convex problem with respect to the variable  $\mathbf{W}^{[n]}$ , and can be evaluated with polynomial computational complexity by the semidefinite programming (SDP) [28], [29] to get the solution  $\hat{\mathbf{W}}^{[n]}$ . However, the solution  $\hat{\mathbf{W}}^{[n]}$  obtained ignoring the constraint (41) becomes in general a full rank matrix, which contradicts the baseline assumption of a rank-1 matrix. To find the rank-1 matrix closest to  $\hat{\mathbf{W}}^{[n]}$ , the rank-1 randomization method [30] is used. Specifically, let us denote the rank of  $\hat{\mathbf{W}}^{[n]}$  by  $r$ , the eigenvalues by  $\lambda_i$ ,  $i = 1, \dots, r$ , with the order  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_r > 0$ , and the corresponding eigenvectors by  $\mathbf{q}_i \in \mathbb{C}^{(N_C \times N_T) \times 1}$ . Then, the rank-1 approximated solution can be written by  $\mathbf{W}^{[n]*} = \lambda_1 \mathbf{q}_1 \mathbf{q}_1^H$ , and the solution for the beamforming vector for given  $\theta$ , is given by

$$\tilde{\mathbf{w}}^{[n]*} = \sqrt{\lambda_1} \mathbf{q}_1. \quad (42)$$

Now, the final step is to find the maximum  $\theta$  that results in a feasible solution of  $\tilde{\mathbf{w}}^{[n]*}$ , which can be readily obtained by a line search algorithm such as the bisection line search.

### Optimization of the UL Power Allocation

As seen from Fig. 3, each BS acquires the UL channel information for all its serving users and exchanges it to all other BSs. The UL PA problem is then formulated and solved at the BS side, and the solution is forwarded to the users.

In the UL time slot, all the users simultaneously transmit UL signals using the energy harvested in the previous time slots. The UL power and corresponding receiver beamforming vector optimization at BS in cell  $j$  for the max min-SINR is formulated by

$$\mathcal{P}_{\text{U1}} : \quad \max_{\mathbf{u}_{j,\pi(j,m)}, p_{j,\pi(j,m)}} \min_{j,\pi(j,m)} \left\{ \rho_{j,\pi(j,m)} \right\} \quad (43)$$

$$\text{s.t. } 0 \leq p_{j,\pi(j,m)} \leq \nu_{j,\pi(j,m)} + \Delta_{j,\pi(j,m)}, \quad (44)$$

$$\left\| \mathbf{u}_{j,\pi(j,m)} \right\|^2 = 1, \quad (45)$$

$$m = 1, 2, j = 1, \dots, N_C. \quad (46)$$

The variables  $\mathbf{u}_{j,\pi(j,m)}$  and  $p_{j,\pi(j,m)}$  in  $\mathcal{P}_{\text{U1}}$  are coupled and need to be jointly optimized, and the problem is non-convex. To make  $\mathcal{P}_{\text{U1}}$  tractable, we find the solution by finding the optimal  $\mathbf{u}_{j,\pi(j,m)}$  for given  $p_{j,\pi(j,m)}$  and then vice versa iteratively. Inserting the optimal beamforming vector known as the MMSE receiver [27] into (22), the UL SINR for given  $p_{j,\pi(j,m)}$  is written by

$$\tilde{\rho}_{j,\pi(j,m)} = p_{j,\pi(j,m)} \left( \mathbf{h}_{j,(j,\pi(j,m))}^{[3]} \right)^H (\mathbf{Z}_{j,\pi(j,m)})^{-1} \mathbf{h}_{j,(j,\pi(j,m))}^{[3]}, \quad (47)$$

and the optimization of  $p_{j,\pi(j,m)}$  can be formulated by

$$\mathcal{P}_{\text{U2}} : \quad \max_{p_{j,\pi(j,m)}} \min_{j,\pi(j,m)} \left\{ \tilde{\rho}_{j,\pi(j,m)} \right\} \quad (48)$$

$$\text{s.t. } 0 \leq p_{j,\pi(j,m)} \leq \nu_{j,\pi(j,m)} + \Delta_{j,\pi(j,m)}, \quad (49)$$

$$j = 1, \dots, N_C, m = 1, 2.$$



Since the cost function (48) is still non-convex, we introduce an additional variable  $\omega$  and modify the cost function as

$$\begin{aligned} \max_{p_{j,\pi(j,m)}, \omega} \quad & \omega \\ \text{s.t.} \quad & \tilde{\rho}_{j,\pi(j,m)} \geq \omega. \end{aligned} \quad (50)$$

Therefore, for given  $\omega$  and with (50), the problem  $\mathcal{P}_{U2}$  can be formulated by

$$\mathcal{P}_{U3}(\omega) : \text{Find } p_{j,\pi(j,m)} \quad (51)$$

$$\text{s.t. } \tilde{\rho}_{j,\pi(j,m)} \geq \omega, \quad (52)$$

$$0 \leq p_{j,\pi(j,m)} \leq \nu_{j,\pi(j,m)} + \Delta_{j,\pi(j,m)}, \quad (53)$$

$$j = 1, \dots, N_C, m = 1, 2.$$

The feasibility problem  $\mathcal{P}_{U3}(\omega)$  is now in a linear problem form, and can be readily solved by LP with polynomial time. The maximum  $\omega$  resulting in a feasible solution of  $\mathcal{P}_{U3}(\omega)$  can be obtained using a linear search such as the bisection method.

After finding the optimal UL power  $p_{j,\pi(j,m)}^*$  for user  $\pi(j,m)$ , each BS announces it to their serving users and the users transmit UL data with the power  $p_{j,\pi(j,m)}^*$ . The unused power is stored at the battery, and the battery level is updated as

$$\Delta_{j,\pi(j,m)} \leftarrow \nu_{j,\pi(j,m)} + \Delta_{j,\pi(j,m)} - p_{j,\pi(j,m)}^*. \quad (54)$$

By using the proposed update algorithm in (38), unnecessary energy consumption at the BSs can be minimized, thereby improving the energy efficiency. Finally, we propose a cascaded DL and UL design as shown in Fig. 5.

## Numerical Results

For comparison of the average minimum user rate, the random beamforming and DL sum-rate maximizing [18] schemes are considered as baseline schemes, in which the beamforming vector is randomly designed and jointly designed across all the cells to maximize the DL sum-rate, respectively. In addition, the maximum EH scheme with and without the minimum SINR constraint [20] are considered. Moreover, we also evaluate the proposed scheme with and without UL PA. It is assumed that the average SNR is the same for all the channels, and that each channel coefficient is an i.i.d. complex Gaussian random variable with zero mean and unit variance. It is also assumed that  $\Lambda_{j,n} = P \in [0, 2]$ ,  $\gamma_{j,n} = 0.5$  and  $\boldsymbol{\pi}_j = [1, 2]$ ,  $j = 1, \dots, N_C$ ,  $n = 1, 2$ .

Fig. 6 demonstrates the boundaries of DL and UL achievable minimum user rate pairs for  $N_C = 3$ ,  $N_T = 2$ , and SNR = 5dB and 25dB. The choice of the parameter  $P$  in the proposed scheme and the choice of the minimum SINR constraint in the max EH with a minimum DL SINR constraint scheme determine respective maximum achievable DL-UL minimum user rates on the rate boundary. The convex hull inside the rate boundary of each scheme is also achievable

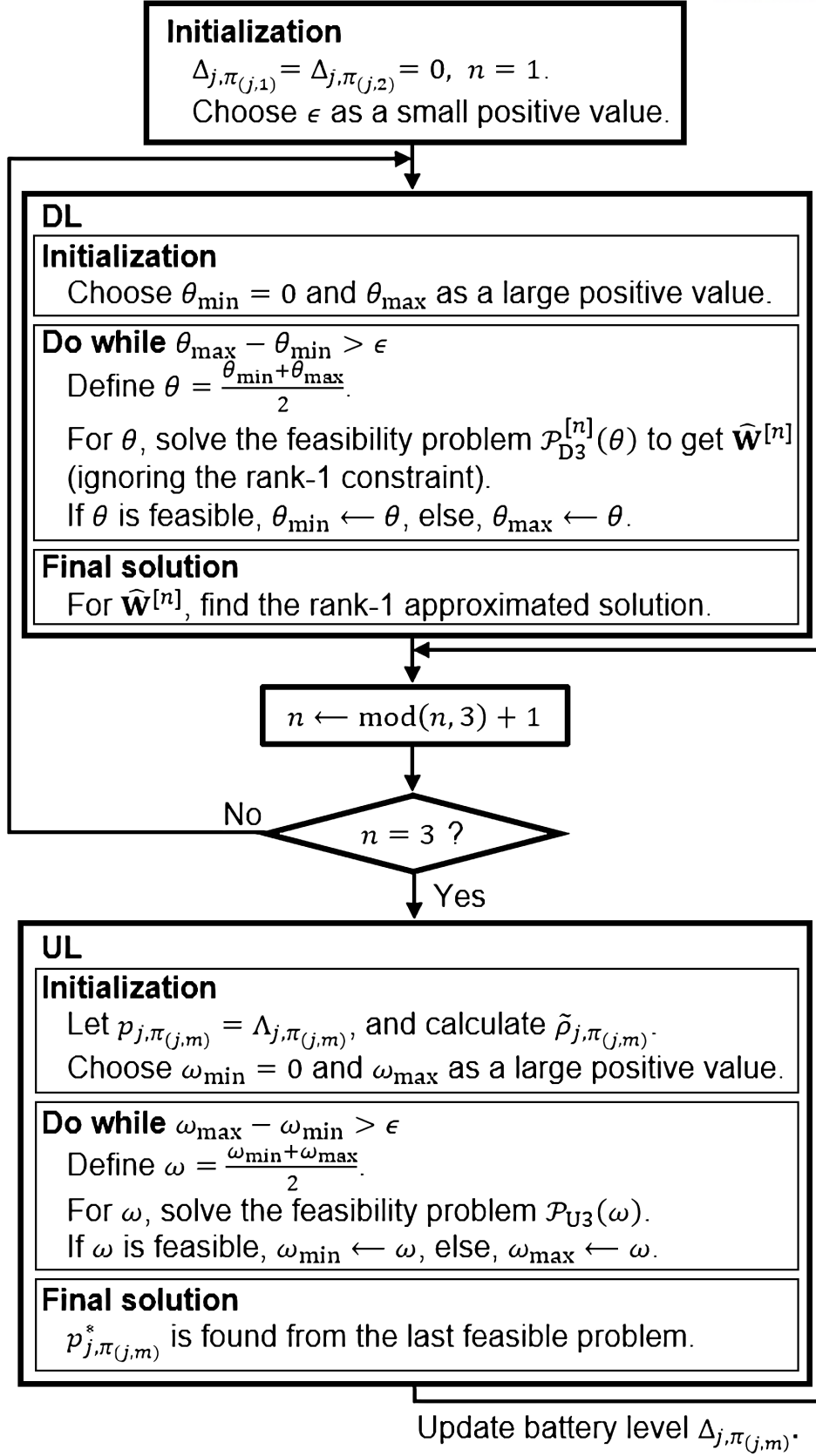


Figure 5: Proposed max-min-SINR DL beamforming and UL PA design

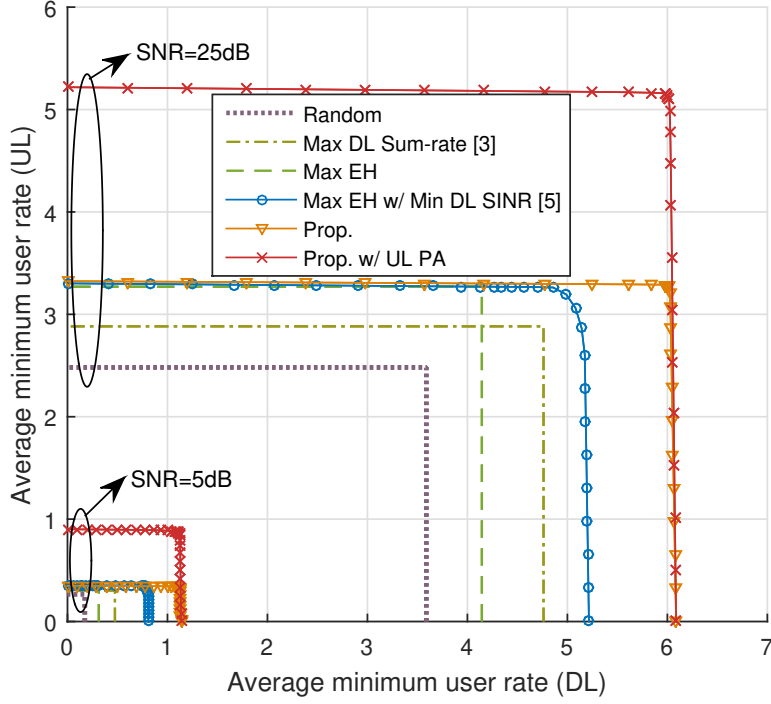


Figure 6: Minimum UL user rate vs. minimum DL user rate

by controlling the time duration of DL and UL. As seen from Fig. 6, the proposed scheme even without UL PA outperforms the existing schemes in terms of the minimum DL and UL rates for all SNR regime, since the DL beamforming design can be more emphasized for maximizing the minimum DL rate owing to the consideration of unused UL power in the problem.

With UL PA, the proposed scheme achieves much broader rate region than the existing schemes, which in turn shows that merely maximizing the harvested energy does not guarantee a high achievable rate due to UL ICI and that the UL PA is essential to significantly improve the UL rate.

Fig 7. shows  $\theta_{\min}$  and  $\theta_{\max}$  versus the number of iterations of the ‘do-while’ part of DL in Fig. 5, i.e., the number of updates for  $\theta$  needed to solve  $P_{D3}^{[n]}(\theta)$ , for  $N_C = 3, N_T = 2$ , and SNR = 5dB, 10dB and 15dB. As shown in Fig. 7, the proposed iterative DL beamforming design converges to the solution within 10 iterations for all SNR regime. It can be also shown by numerical simulations that the proposed iterative UL PA converges within a few iterations.

Table 1 shows the relative energy efficiency defined as the ratio of the minimum UL user rate to the total UL power consumption in comparison to that of the maximum EH scheme with the minimum SINR constraint. Due to the cascaded DL and UL design with the consideration of the unused UL power, UL power can be significantly saved in the proposed scheme achieving the same or even higher minimum UL rate as seen from Fig. 6. As a result, the proposed scheme with UL PA shows the highest UL energy efficiency among all compared schemes for all SNR regime, as seen from Table 1.

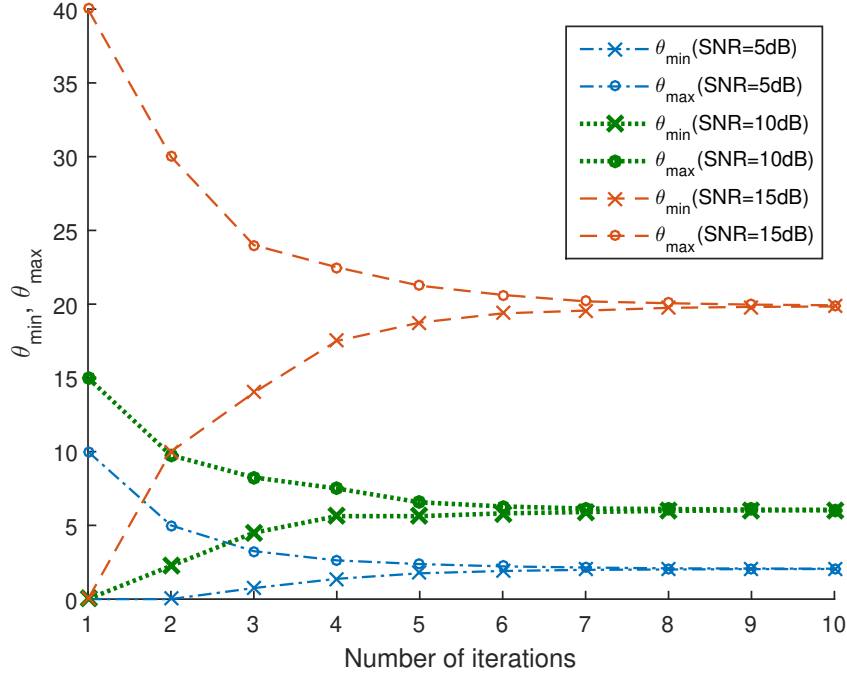


Figure 7:  $\theta_{\min}, \theta_{\max}$  vs. number of iterations

Table 1: UL energy efficiency

	5dB[%]	25dB[%]
Random	93.5	93.3
Max DL Sum-rate [18]	94.1	98.3
Max EH	124.7	135.1
Max EH w/ Min DL SINR [20]	100	100
Prop.	143.3	162.1
Prop. w/ UL PA	173.0	206.4

## Conclusion

The general communication and energy transfer scheme for multicell networks composed of BSs with multiple antennas and users each with a single antenna has been proposed. In the proposed scheme, the concepts of SWIPT and WPCN are jointly considered for the DL and UL time slots, and the UL PA is also considered in the UL time slots. Simulation results show that the proposed scheme not only achieves larger minimum DL and UL SINRs region but also exhibits much improved energy efficiency.

## 2.3 RNN-Based Node Selection for Sensor Networks with Energy Harvesting

### Introduction

Recently, as smart home and smart factory have emerged, sensor networks consisting of devices with energy harvesting circuits have attracted worldwide attention. Efficient energy management makes the network sustainable and environment-friendly. In [31], the authors consider the sensing utility maximization problem where each node transmits sensing data with harvested energy from ambient environment. In [32], the authors solve the minimization problem of the packet delay where a transmitter harvests energy via wireless power transfer. However, these previous schemes merely focus on either uplink (UL) or downlink (DL) data transmission. In fact, the unpredictability of future channel condition and the coupled DL and UL design make joint optimization problems more challenging. There is an approach to optimal design for both simultaneous wireless information and power transfer (SWIPT) and wireless powered communication network (WPCN) in cellular networks [33]. The major problem in sensor networks is UL/DL node selection, which becomes more difficult to solve as the number of nodes increases enormously in the near future.

In this paper, we address UL/DL node selection problem with limited information including future channel state and each node's situation such as battery levels and packet deadlines. We then propose a recurrent neural network (RNN) based algorithm for the node selection. To the best of our knowledge, this is the first work that jointly considers DL SWIPT and UL WPCN for sensor networks. Simulation results show that the proposed scheme outperforms other existing schemes.

The remainder of this paper is organized as follows. Section II explains the background, and Section III describes the system model. Section IV presents a method to predict the future decision using RNN. Section V provides simulation results, and Section VI concludes the paper.

### Background

- **IEEE 802.15.4g**

The IEEE 802.15.4 [34] is the standard for Low-Rate Wireless Networks. One important variant is the IEEE 802.15.4g [35] Smart Utility Network (SUN) supports various PHY layers: frequency shift keying (FSK) PHY, orthogonal frequency division multiplexing (OFDM) PHY, and offset quadrature phase-shift keying (O-QPSK) PHY. In this paper, we assume PHY layers are OFDM PHY.

Active Superframe Portion					Inactive Portion
Beacon	CAP	EB	CAP	CFP	

Figure 8: Superframe structure

As seen in Fig.8, an active portion of a superframe consists of contention access periods (CAPs) and a contention free period (CFP). For low-latency applications, a master node (MN) in the SUN dedicates CFP to those applications including itself. This paper considers which node will be selected as the user of the CFP by the MN.

### • Recurrent Neural Network

RNNs are artificial neural networks that recognize patterns in sequence or time series data

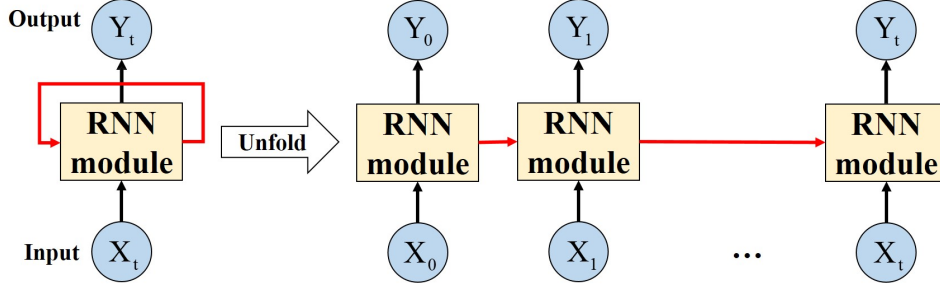


Figure 9: Basic structure of RNN

in the form of texts, genes, handwriting, voice signals, sensor-detected data, stock quotes, etc. Fig.9 presents a basic structure of RNN.  $X_t$  and  $Y_t$  are the input node which receives data from outside of the network and the output node which yields results respectively. The RNN module consists of hidden states. Since the RNN is a directed cycle artificial neural network with hidden states connected with the directional edges (red arrows in Fig.9), past outcomes can affect future outcomes.

Depending on the number of inputs and outputs, the RNN can be formed in a variety of structures, such as one-to-many, many-to-one, and many-to-many. In this paper, since we know the previous and the current channel information and find out what a decision the MN should make, we follow the many-to-one RNN structure.

### System model

In this paper, we consider a sensor network consisting of a MN and  $N_{SN}$  slave nodes (SNs) as shown in Fig.10. The MN has a reliable power and each SN is powered by energy harvested from the ambient RF signals and the ambient environment. We assume that each SN has a battery to store the harvested energy, and it is consumed by sensing data and transmitting them to the MN. The MN always has packets to transmit to each SN, and each SN has packets to transmit to the MN. Every packet has its own deadline and the each deadline follows a distribution. We consider finite time slotted system. For each time slot, an UL or a DL may occur, and only one node who is selected by the MN can transmit or receives signals. While a selected SN transmits UL data by using the energy harvested or receives DL data, other SNs can harvest ambient RF signals or environment. The channel coefficients between the MN and SN  $i$  at the time slot  $t$  is denoted by  $h_i^{[t]} \in \mathbb{C}$ ,  $i = 1, 2, \dots, N_{SN}$ , and  $t = 1, 2, \dots, T$ . It is assumed that

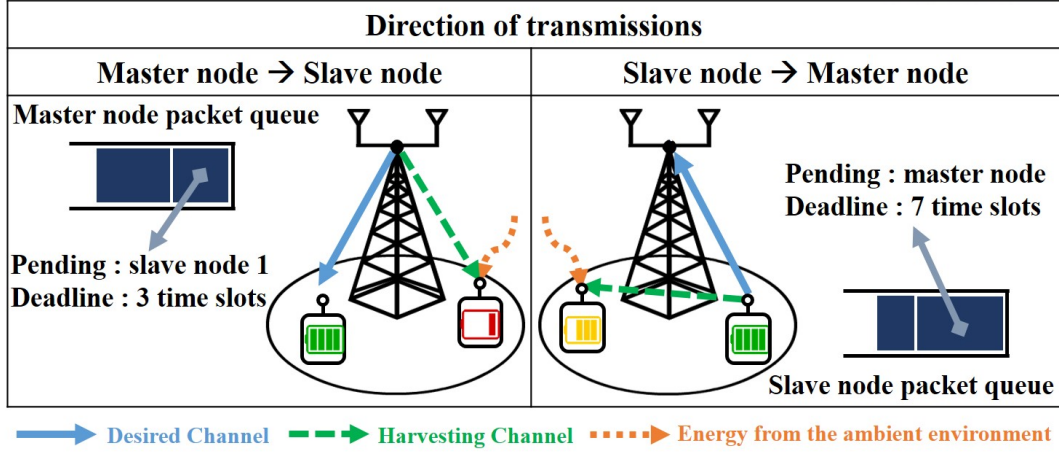


Figure 10: System model

the channel coefficients remain constant within a time slot and then change to another values with correlation at the next time slot, i.e., quasi-static fading. The MN is assumed to be able to acquire its outgoing channels through channel sounding at each time slot, but does not know its future channels.

When SN  $i$  is selected to receive DL data, the received signal is written by

$$y_i^{[t]} = \sqrt{P_m} h_i^{[t]} x_i^{[t]} + z_i, \quad (55)$$

where  $P_m$  is the transmit power of the MN,  $x_i^{[t]}$  is the unit-variance transmit symbol at the time slot  $t$ , and  $z_i$  is the additive white Gaussian noise (AWGN) at SN  $i$  with zero mean and variance of  $N_0$ . If the transmitted DL data is already missed deadline, we impose a DL deadline penalty. Unselected SNs harvest energy and the amount of harvested energy is written by

$$\gamma P_m h_{i-1}^{[t]*} h_{i-1}^{[t]} + P_e, \quad (56)$$

where  $i_{-1}$  denotes an index which means all SNs except for SN  $i$ .  $0 < \gamma \leq 1$  denotes the harvesting efficiency and  $P_e$  is energy randomly harvested from the ambient environment.

When SN  $i$  is selected to transmit UL data, if the SN has enough battery power, it transmits UL data with the battery level and UL packet deadline ((i) in Fig.11), otherwise it transmits battery level and UL packet deadline only with limited bits ((ii) in Fig.11). If the battery power is not enough or the UL deadline is missed, we impose a battery level penalty or an UL deadline penalty. Unselected SNs harvest energy from the transmitted signal by SN  $i$  and the energy from the ambient environment. Since the MN doesn't know the channel coefficients between SNs, we assume that the amount of energy harvested for each unselected SN is randomly set within a reasonable range.

Finally, the objective of this system is to minimize the summation of the battery level and UL/DL deadline penalty.

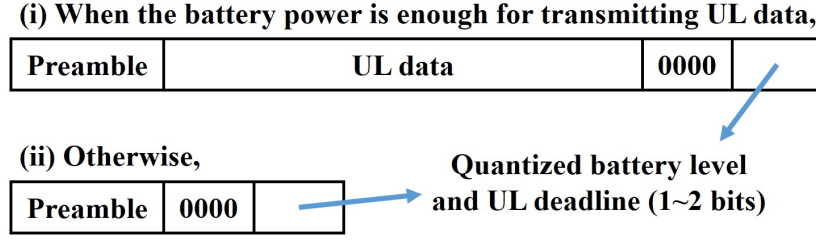


Figure 11: UL packet format

### Approach: predict the best decision

The goal of RNN-based algorithm is to guess the best decision combination that minimizes the number of penalties for  $T$  time slots, when only the previous and the current channel information are given. Each decision for each time slot can be UL/DL of SN  $i$ . The number of cases for all decision combinations is  $N_{DC} = (2 \times N_{SN})^T$ . Since the MN does not know each battery level and packet deadline of each SN, we do not take as input those things. Fig.12 shows an example of the RNN structure of the proposed scheme when  $T = 3$  and  $N_{SN} = 2$ .

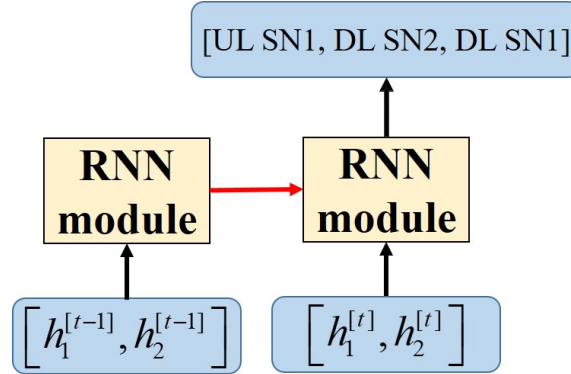


Figure 12: RNN structure of the proposed scheme

To make labeled ground-truth, first of all, we generate enormous number of  $T$  time slots channel coefficients and multiply it by the conjugate form as (57).

$$\begin{bmatrix} h_1^{[1]*} \times h_1^{[1]} & h_1^{[2]*} \times h_1^{[2]} & \dots & h_1^{[T]*} \times h_1^{[T]} \\ \vdots & \vdots & \ddots & \vdots \\ h_{N_{SN}}^{[1]*} \times h_{N_{SN}}^{[1]} & h_{N_{SN}}^{[2]*} \times h_{N_{SN}}^{[2]} & \dots & h_{N_{SN}}^{[T]*} \times h_{N_{SN}}^{[T]} \end{bmatrix} \quad (57)$$

For one channel realization, we assume the various initial battery levels of SNs and UL/DL packet deadline conditions, and then look for the best decision considering all cases. Each channel realization has a best decision combination for  $T$  time slots. After generating lots of channel realizations and finding the corresponding the best decisions, we quantize the elements of (57) as {low, middle, high} to simplify the cases. Then, (i) classify the same first  $(T - 1)$  time slots of (57) with the same value (green dashed box in Fig.13). (ii) Find the most frequent



occurrence in the last time slot of it to predict the most likely future channel conditions (red box in Fig.13). Once the most likely future channel conditions have been determined, (iii) the most frequent decision among the corresponding decisions can be found (green box in Fig.13).

Let us define a ground-truth class  $\mathbf{U} = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_{N_{DC}}]$  as the one hot coded  $N_{DC}$  decision combinations. For example,  $\mathbf{u}_1 = [1, 0, \dots, 0]$  is the first decision combination. Finally, decision labeled ground-truth for the quantized previous and current channel information is completed as (iv) in Fig.13.

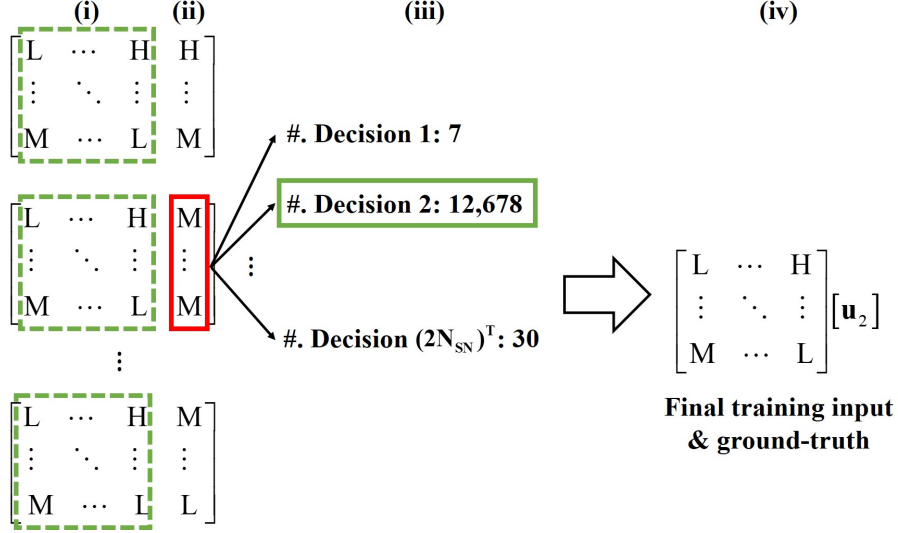


Figure 13: Method to make labeled ground truth

When training the RNN, we take inputs as the quantized previous and current channel information. The network outputs a probability distribution,  $\mathbf{p} = (p_1, p_2, \dots, p_{N_{DC}})$ , over  $N_{DC}$  cases, that is the number of all decision combinations. The probability distribution  $p$  is computed by a softmax over the  $N_{DC}$  outputs of a fully connected layer.

Each training input is labeled with a ground-truth  $\mathbf{u}_d$ . We use a mean square error loss  $L$  on each labeled training input to train:

$$L(\mathbf{p}, \mathbf{u}_d) = \frac{1}{N_{DC}} \sum_{k=1}^{N_{DC}} (p_k - u_{dk})^2 \quad (58)$$

in which the  $p_k$  is the  $k$ -th element of the  $\mathbf{p}$  and the  $u_{dk}$  is the  $k$ -th element of the  $\mathbf{u}_d$ .

## Simulation results

With the saved weights from the training, we can test our network. Unlike training, we take input as unquantized previous and current channel information. It is assumed that the average SNR is the same for all the channels, and that each channel coefficient is an i.i.d. complex Gaussian random variable with zero mean and unit variance. It is also assumed that  $N_{SN} = 2$ ,  $T = 3$ ,  $\gamma = 0.7$ ,  $P_m = 3\text{mW}$ , and  $P_e = 0.5\text{mW}$ . The number of decision combinations  $N_{DC}$  is automatically 64.

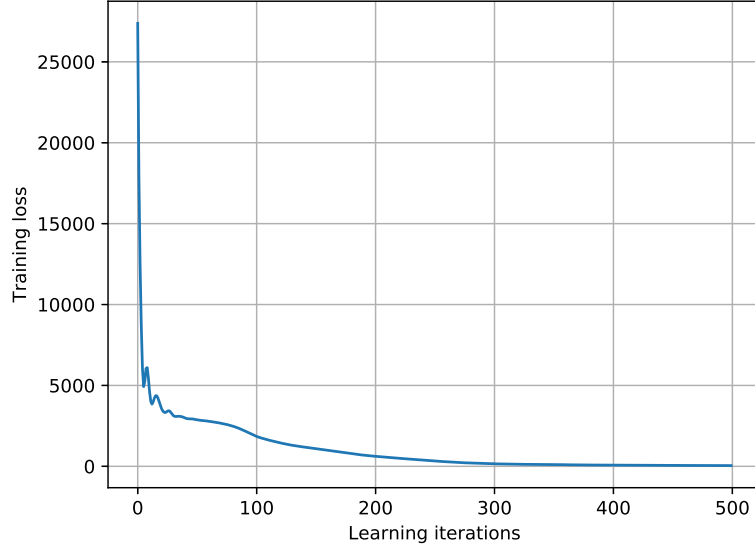


Figure 14: Training loss vs. Learning iterations

We use training optimizer as adaptive moment estimation (Adam) optimizer. Step size is 0.01 and two exponential decay rates for the moment estimates are 0.9 and 0.999. The number of training data is  $2 \times 10^4$  and the number of iterations for the training is also  $2 \times 10^4$ . Fig.14 presents the training loss versus learning iterations. It shows the training loss decreases in tens of times of learning iterations. After  $3 \times 10^2$  iterations, the training loss does not change much and converges to almost zero. If the channel conditions are similar to those of the training input, it can be expected that we can get decision results like labeled ground-truth with high probability.

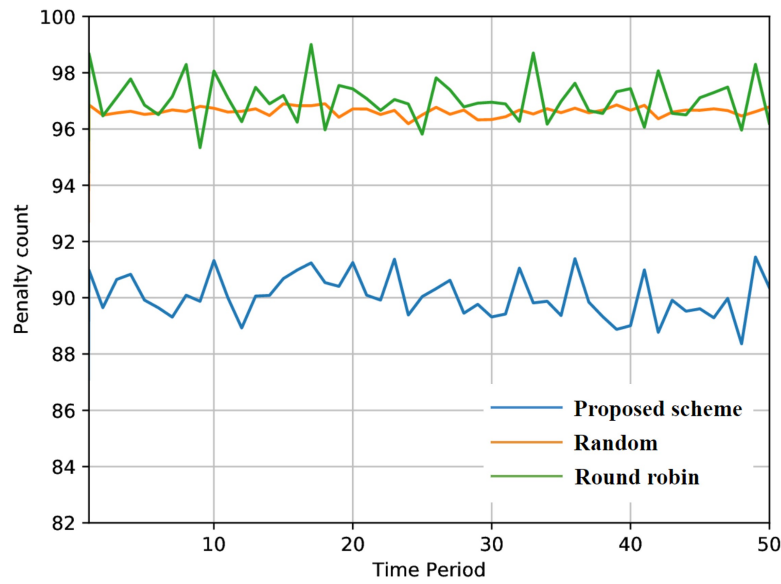


Figure 15: Number of penalties vs. time

Fig.15 shows the averaged penalty count. We perform  $10^3$  experiments, in which the number of time period is  $5 \times 10^3$ . Each penalty count is measured every  $10^2$  time period and then averaged over  $10^3$  experiments. For comparison of the penalty count, the random selection and the round robin selection are considered as baseline schemes, in which the final decision is randomly selected and taken in turns all decisions without a priority. As seen in Fig.15, the proposed scheme outperforms the base schemes in terms of the penalty count for all time period, since the proposed scheme predicts future channel implicitly and tries to determines the best decision for the future.

## Conclusion

We have proposed a new RNN-based node selection scheme that selects a node that transmits UL data or receives DL data in the future time slot, where the decision is determined toward minimizing penalty count. The determination of these decisions also takes into account battery level, UL/DL deadline as well as future channel information implicitly. Through simulation results, we have shown that the proposed RNN-based node selection scheme outperforms the other base schemes.

### III Learning-based Robot Vision Systems

#### 3.1 Privacy-Preserving Robot Vision with Anonymized Faces by Extreme Low Resolution

##### Introduction

Robot camera systems are becoming more important as a component of robot perception for autonomous navigation, robot-robot or robot-human interactions, surveillance, and etc. As small-scale mobile robots, such as drones, planes, and rovers are expected to be pervasive to meet increasing demands of new customer services such as unmanned delivery, 24/7 surveillance, internet connectivity in disaster, and etc, there is an increasing concern of privacy threat from recording all unwanted images from the cameras on robots. Thus, human privacy protection from robot cameras is a serious and real social challenge [36], while we do not want to sacrifice robot perception performance. Indeed, several privacy-attacks have been reported, such as IP cameras cracked by human hackers and private images from home cameras leaked through the internet.

In this work, we pursue two goals for robot camera systems: (1) no one can access any of privacy-sensitive visual information (e.g., faces) by cracking or installing backdoors for fundamental privacy protection, and (2) robots can benefit from the video for their perception as much as they can without any consideration of privacy protection. Specifically, we develop a mobile robot system with a novel feature of privacy-preserving face detection. Unlike the previous approaches that blur or anonymize the face blocks from high resolution (HR) images, face blocks are detected from extreme low resolution (LR) images using the proposed deep learning-based algorithm. Simultaneously, to improve the performance of the robot perception, the resolution of the privacy-insensitive blocks of images is allowed to increase. The proposed face detection system guarantees that all the face blocks in an image are designed to never be recognizable in any of processing or memory, thus providing fundamental privacy protection. Furthermore, we propose a pixel-by-pixel post-processing algorithm to distinguish the face and background blocks at a pixel level more accurately, thus ensuring that robots can perform the SLAM as they can without the consideration of face anonymization. We empirically confirm that the proposed extreme LR face detection algorithm outperforms the state-of-the-art technique, and that robots with the proposed face detection system still can perform well ORB-SLAM2.

##### Related work

- **Privacy Protection from Cameras**

A variety of studies have been conducted to meet the social needs of privacy protection from cameras. The work [37] studied scene recognition from images captured with first-person cameras, detecting locations where the privacy needs to be protected. This will allow the device to be automatically turned off at privacy-sensitive locations. There is

another approach for privacy preserving action detection through learning a video face anonymizer [38]. The video anonymizer performs pixel-level modifications to anonymize each person’s face, with minimal effect on action detection performance. However, all these aforementioned methods are based on software-level processing using HR videos. For this reason, they are still not safe from cyber attacks. To deal with this problem, one of the fundamental solutions is taking extreme LR videos without having privacy-sensitive visual information in any memory even during interim processing. In [39–41], the authors studied activity recognition from extreme LR (e.g.,  $12 \times 16$ ) anonymized videos. The work [39] introduced the concept of learning the optimal set of image transformations to generate multiple LR training videos optimized for the activity classification from one HR video. In addition, [40] presented an extreme LR activity recognition model using a two-stream multi-Siamese convolutional neural network. The model made good use of the inherent characteristic of LR videos: Two LR images originated from the exact same scene often have totally different pixel values depending on their transformations. Furthermore, [41] explicitly learned a degradation transform for the original HR video, in order to optimize the trade-off between visual recognition performance and the associated privacy budgets. Although these works are effective for preserving privacy, they sacrifice too much visual information by having all the parts of images at extreme LR. If robots are equipped with cameras taking only extreme LR images or video, then they cannot operate basic applications based on robot vision such as autonomous navigation, interaction and etc.

- **Small Face Detection**

Face detection is a well-studied research field, however, small face detection is still a challenging area. One of the latest works [42] described a detector that can find tiny faces, exploring the role of the contextual reasoning and scale in a pre-trained deep networks. However, it was designed to find faces well only in pre-trained object scales. Another approach to find small faces is to directly generate a clear HR face from a blurry small one by adopting a generative adversarial network (GAN) and then detect the face [43]. In [44], a model that finds a small frontal face is applied to a robot. Since all of these studies find small faces in a HR image of at least  $240 \times 320$ , it is difficult to apply them to finding faces in extreme LR images. Finding a face in extreme LR images is practically difficult because it lacks the amount of information. Even non-visual information such as optical flow [45] and point trajectory [46] available in action recognition is not available for face detection.

- **SLAM**

Simultaneous localization and mapping (SLAM) techniques build a map of an unknown environment and localize the sensor, e.g., a robot or a camera, in the map with a strong focus on real-time operation. In our system, we used ORB-SLAM2 [47], one of the most popular algorithms of the visual SLAM, which mainly uses cameras.

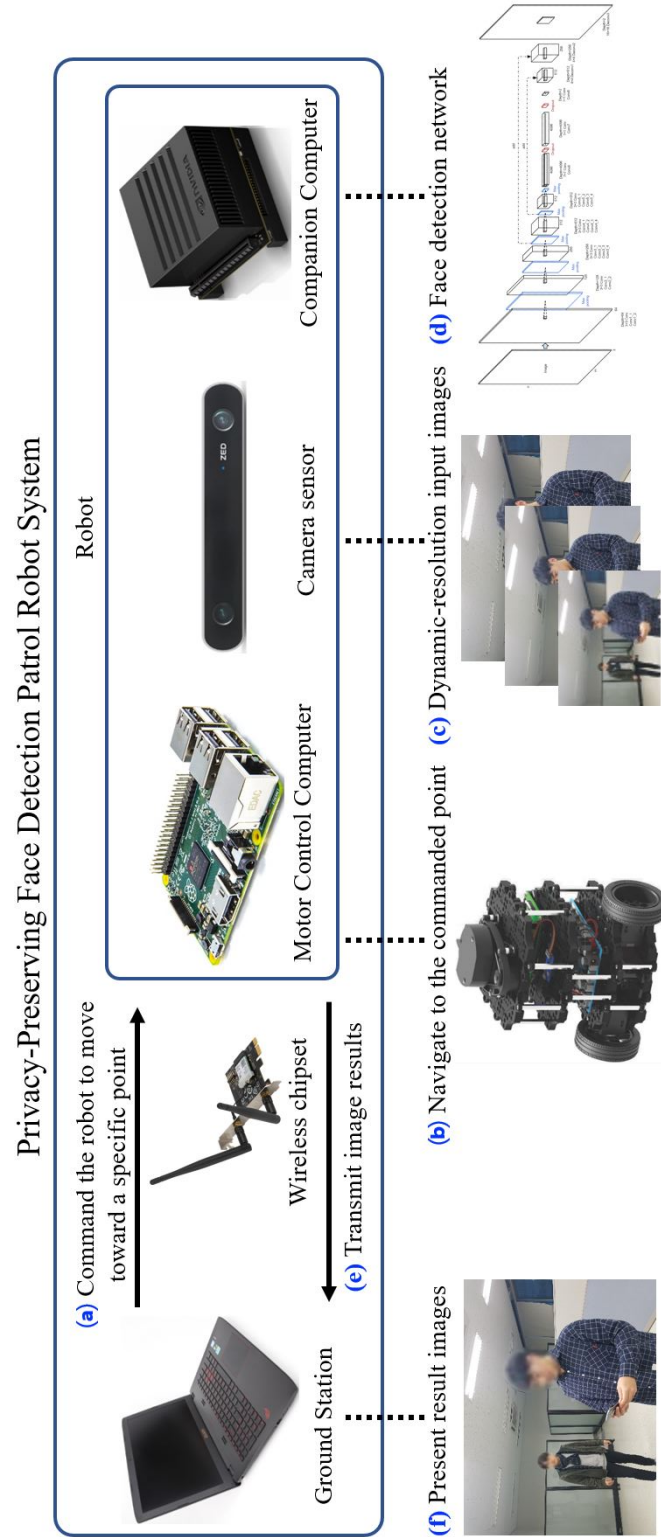


Figure 16: Composition of the developed patrol robot system with privacy preserving face detection.

## Privacy-preserving face detection system

We develop a patrol robot system with cameras continuously recording video for its navigation. Fig. 16 shows the overall composition of our proposed patrol robot system. The ground station wirelessly commands the robot to move toward a specific point. Then, the robot navigates to the commanded point by the motor control computer mounted on it. At the same time, the robot obtains dynamic-resolution input images through the camera module sensor. Then, the companion computer mounted on the robot performs the privacy-preserving face detection. The robot performs ORB-SLAM2 and transmits the resultant image after the face detection back to the ground station wirelessly, and the ground station presents the results to the window screen.

The goal is three-fold:

- Every face in images should never be identifiable by virtue of the proposed privacy-preserving face detection algorithm. Specifically, the resolution of any face block should never be higher than  $15 \times 15$ .
- The robot SLAM should work well as it does without the proposed face anonymization. To this end, the background blocks of images should not be in LR.
- All the computation including the proposed privacy-preserving face detection and SLAM should run on a real-time basis only on the companion computers. This makes our robot system operate fully autonomously even without the connection to the ground station.

Each component is described with details in the sequel.

### A. Face Detection Approach

#### • Limitations of the State-of-the-Art

The face detection in our system is to detect faces from extreme LR images. Though not for privacy protection, several studies have been conducted for small face detection at a distance. The state-of-the-art small face detection [42] scales a single HR snapshot up and down to create multiple snapshots to find tiny faces. This approach employs a set of bounding box shapes which are selected for particular object sizes. However, this method only works well for preset resolutions and bounding box shapes, which cannot be dynamically adopted. As a result, unless the bounding box shapes are various enough to cover all possible face sizes, which is not possible, the bounding box may not cover some pixels of the face, or may cover too many pixels out of the face. The former case leads to a risk of privacy invasion, and the latter case results in a significant loss of information needed for SLAM. For this reason, we develop a framework for finding extreme LR faces more reliably on a pixel-by-pixel basis. In addition, via dynamic resolution control and pixel-by-pixel post-processing, we attempt to increase the resolution of privacy-insensitive background blocks for accurate SLAM.



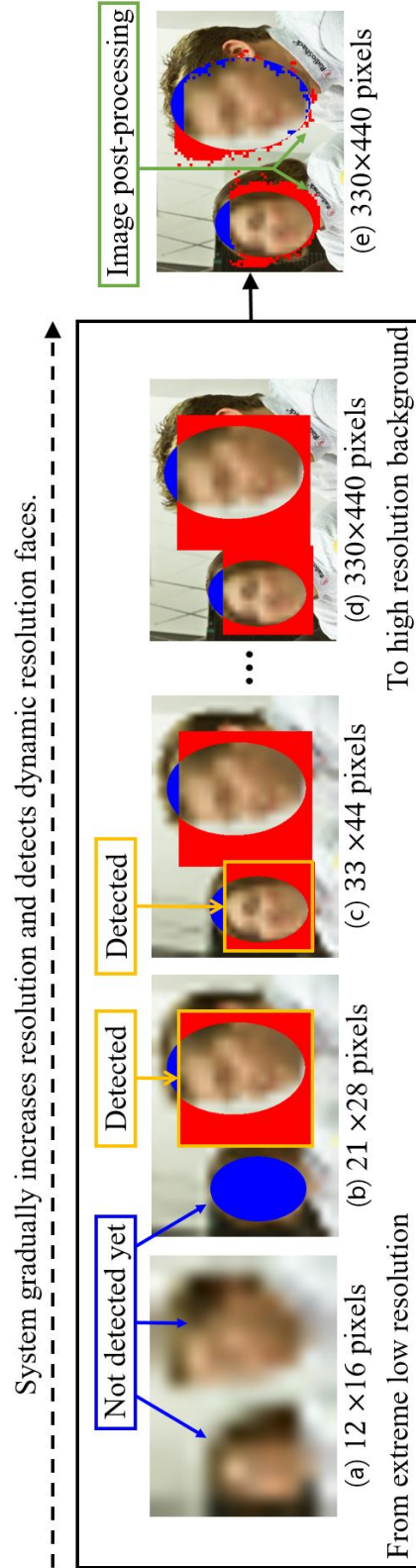


Figure 17: Dynamic resolution face detection architecture.



### • Dynamic Resolution Face Detection

The key idea is every face block should be detected at a resolution lower than  $15 \times 15$ , such that it is not identifiable. The proposed dynamic resolution face detection architecture is illustrated in Fig. 17. When the camera sensor receives an image, the system does not store the images in any computer memory, but directly converts the images to extreme LR, e.g.,  $12 \times 16$ , as in Fig. 17 (a). This conversion can be easily done by dedicated up/down-scaling DSP chips or FPGA circuits attached on the camera sensors. The companion computer mounted on the robot determines if there are faces in the image. If a face is found as shown in Fig. 17 (b), the resolution of that block is maintained as extreme LR in the next step, and the resolution of the rest is slightly increased. In Fig. 17 (b), only the larger face, but still in very LR, is detected, and the smaller face cannot be detected. As Fig. 17 (c) shows, the resolution of the image excluding the larger face block is higher, so that the smaller face, which is still in very LR, now can be found. This procedure is repeated until all the faces are detected in very LR. In the end, only the non-facial parts are in relatively HR as in Fig. 17 (d). Finally, before outputting the results, the system checks again pixel-by-pixel for the parts that have been judged to be faces at LR. If the probability of it not being a face pixel is very high, the system performs post-processing to increase the resolution as shown in Fig. 17 (e).

The aim is to detect faces at the lowest possible resolution to protect privacy. To this end, the system increases resolution only very gradually from an extreme LR and uses the training data that is well tailored for finding faces at very LRs, as described in the sequel.

### • Training Dataset

For training, we selected the AFLW dataset [48], a public dataset popularly used for face detection. Unlike the existing work [42], we do not use the original HR dataset for training, but resize it so that the face parts are of the target size  $15 \times 15$ . The dataset is composed of annotated face images with a large variety in appearance (e.g., pose, expression, ethnicity, age, and gender) as well as environmental conditions, containing about 25k faces.

Resizing an image to make the face within it  $15 \times 15$ -sized is not difficult if there is only one face in the image. However, if there are several faces in one image, a careful consideration should be given when generating the training data. We describe our training data generation process for the example of the image with two faces in Fig. 18. The original image, Fig. 18 (a), is resized so that the resolution of the larger face is  $15 \times 15$  as in Fig. 18 (c). In the resized image, the smaller face is smaller than  $15 \times 15$  pixels as in Fig. 18 (c), and labeled as “face” without any problem. This resized image is used as training data. Note that finding a face smaller than  $15 \times 15$  pixels only makes the face more unidentifiable. To make the smaller face in the original image detectable at the solution of  $15 \times 15$ , the original image is additionally resized so that the smaller face is  $15 \times 15$ -sized as in Fig. 18 (b). This resized image now includes the larger face with size larger than  $15 \times 15$  as in Fig.

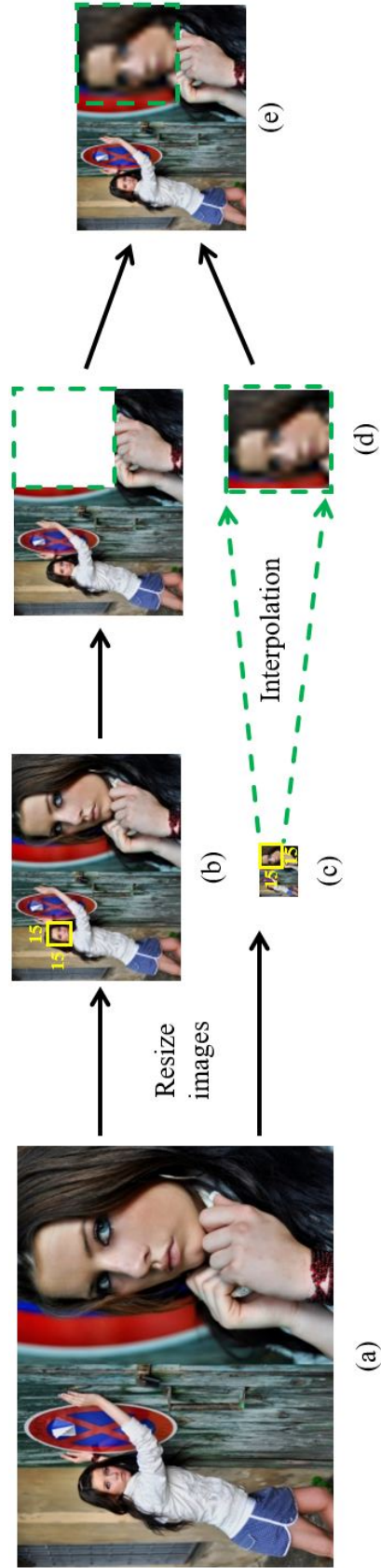


Figure 18: Our training data generation process for the example of the image with two faces.

18 (b), which may result in a privacy threat. To handle this problem, the  $15 \times 15$ -sized larger face part of Fig. 18 (c) is upscaled by the bicubic interpolation so that the upscaled image replaces the larger face part of Fig. 18 (b) as in Fig. 18 (d). Finally, additional training data is generated as in Fig. 18 (e), where all the faces are  $15 \times 15$ -sized. Note that training data generation is performed offline using the public dataset, and hence there is no privacy issue.

- **Deep Learning**

We use a fully convolutional networks (FCN) [49] structure which is one of the effective neural networks for semantic segmentation. Its feature extractor is based on VGG19 [50]. FCN predicts dense outputs from arbitrarily-sized inputs; that is, in-network upsampling layers enable pixelwise prediction, and the network handles inputs of various sizes. To improve the performance of these features, the following tasks are performed.

- **Bounding Box**

Due to the nature of the FCN structure, when the system makes decisions at LR on a pixel-by-pixel basis, there is a possibility that only a few pixels in the face block are judged to be facial while many pixels in the same face block are judged to be non-facial, which results in a privacy threat. For instance, a donut-like face shape can be detected, which makes no sense. To avoid this limitation, the system calculates the probability of being facial for each pixel and then re-calculates an averaged probability value for the group of pixels within a bounding box sliding throughout the whole image. If the averaged value exceeds a threshold, the bound box is judged to be a face. We have empirically confirmed that a fixed size bounding box of  $15 \times 15$  pixels gives us reliable results, whereas the existing scheme [42] employs a variety of bounding box shapes and sizes.

- **Image Post-processing (IPP)**

Using a bounding box provides reliable face pixel detection, but sometimes many non-facial pixels are included in the bounding box, causing many false positives. In order to restore the resolution of the non-facial pixels that have been judged incorrectly, we propose the system checks again the parts that were judged to be faces, just before it outputs the final image. It increases the resolution of the pixels with the probability of not being a face exceeds a threshold, e.g., 99%, as in Fig. 17 (e).

- **Challenge: Real-time Processing**

The existing approach [42] basically uses HR images for both training and testing. In addition, the method scales images up and down to create multiple templates, finds faces, and then combines the results. Therefore, the computational complexity is relatively high and the run-time is 1.4FPS on 1080p resolution and 3.1FPS on 720p resolution based on the Resnet101-aided detector. On the other hand, our network has relatively low computational complexity since it starts from extreme LR. Though the proposed system

increases resolution gradually, there is still a benefit in computational complexity because at any step it does not reprocess the face parts previously found. When the resolution gradually increased 6 times, from  $12 \times 16$  to the original full resolution, the run-time is 6.996 FPS. Furthermore, in contrast to proposal-based detectors such as Faster R-CNN [51], which employs as many bounding boxes as the number of proposals, our run-time does not increase with the number of faces in an image.

## B. Robot Control

The robot system consists of the following components: (1) **ROS-based motor control computer** operating the motors so that the robot can move physically; (2) **wireless chipset** receiving commands from the ground station and transmitting image results to the ground station; (3) **camera sensor module** receiving input images to measure the distance between the robot and obstacles when the robot performs SLAM; (4) **companion computer** running the proposed face detection neural network and performing SLAM.

The robot is commanded by the ground station to move to a specific point via the wireless chipset. According to the received commands, the ROS-based motor control computer controls the robot to drive. While the robot is moving, images are continuously received from the camera, which are passed through the face detection neural network at the companion computer on a real-time basis. The privacy-preserved images are sent back to the ground station via the wireless chipset so that we can immediately check the results on the screen of the ground station. At the same time, the robot executes the ORB-SLAM2 algorithm using the processed images.

## Experiments

### A. Face Detection Deep Learning Evaluation

- **Test Dataset**

We tested our model on the AFLW and FDDB [52] dataset. In order to gradually increase the resolution of the images in testing, we had to resize the images with various resolutions in advance. When resizing the images, we could not resize all the images to the same size such as  $24 \times 32$ , because the aspect ratios of the images are not constant. Instead of a fixed aspect ratio, the images were resized based on the sum of the horizontal and vertical pixels such as  $w+h=50$  pixels and  $w+h=75$  pixels.

- **Existing Scheme**

The approach presented in [42] is to utilize multiple templates with the coarse image pyramid, where one template is tuned for 40-140px tall faces and the other one is tuned for less than 20px tall faces. The neural network trained in [42] is implemented for comparison of the baseline scheme. For fair comparison, the same dynamic resolution control is applied to the baseline scheme as in the proposed face detection algorithm, where the resolution of currently non-facial pixels is gradually increased to detect all small faces in the image.

Table 2: Performances of different methods

Dataset	Method	DETECTED	FP	DETECTED-FP
AFLW	Hu et al.	0.688	0.063	0.625
	Proposed w/o IPP	0.922	0.247	0.675
	Proposed w/ IPP	0.896	0.183	0.713
FDDB	Hu et al.	0.479	0.078	0.401
	Proposed w/o IPP	0.643	0.242	0.401
	Proposed w/ IPP	0.608	0.141	0.467

### • Face Detection Results

Taking careful consideration on the privacy problem when detecting faces, we should declare that a face is not well detected even though it is detected in the resolution higher than  $15 \times 15$ . Therefore, instead of using the term IOU (intersection over union) which is a widely used metric, we define four metrics for performance evaluation of the face detection algorithms as follows. (1) DETECTED: The ratio of the number of pixels judged as a part of a face when the size of the face is smaller than or equal to  $15 \times 15$  pixels to the number of pixels actually included in the face. (2) MISSED: The ratio of the number of pixels judged as a part of a face when the size of the face is larger than  $15 \times 15$  pixels, or not judged as a part of any faces until the resolution of the image is fully raised to the number of pixels actually included in the face. (3) FALSE POSITIVE (FP): The ratio of the number of pixels judged as a part of a face at any resolution to the number of pixels actually not included in the face. (4) TRUE NEGATIVE (TN): The ratio of the number of pixels judged as not a part of faces at all resolutions to the number of pixels actually not included in the face.

Table 2 shows the performance of different methods tested with the AFLW and FDDB dataset. It is observed that our methods both with and without IPP show significantly higher DETECTED performance than the approach presented in [42] in both datasets, although they have slightly increased FP values. From the privacy protection perspective, it is much more important than anything else that as many face pixels are found in extreme LR as possible. Thus, our method is more suitable for privacy-preserved robot vision.

Fig. 19 shows the results comparing our proposed method with the approach presented in [42] for three different example images. From the left-hand side column, the images in the figure are the original images, the results of the approach presented in [42], and the results of the proposed face detection algorithm, respectively. In each image, the LR parts indicate DETECTED, and the red and blue parts indicate FP and MISSED, respectively. Two numbers below each image represent DETECTED and FP in percentage. The proposed system finds faces with bounding boxes, gradually increases the resolution, and finally increases the resolution of the pixels with high probability of being non-facial.



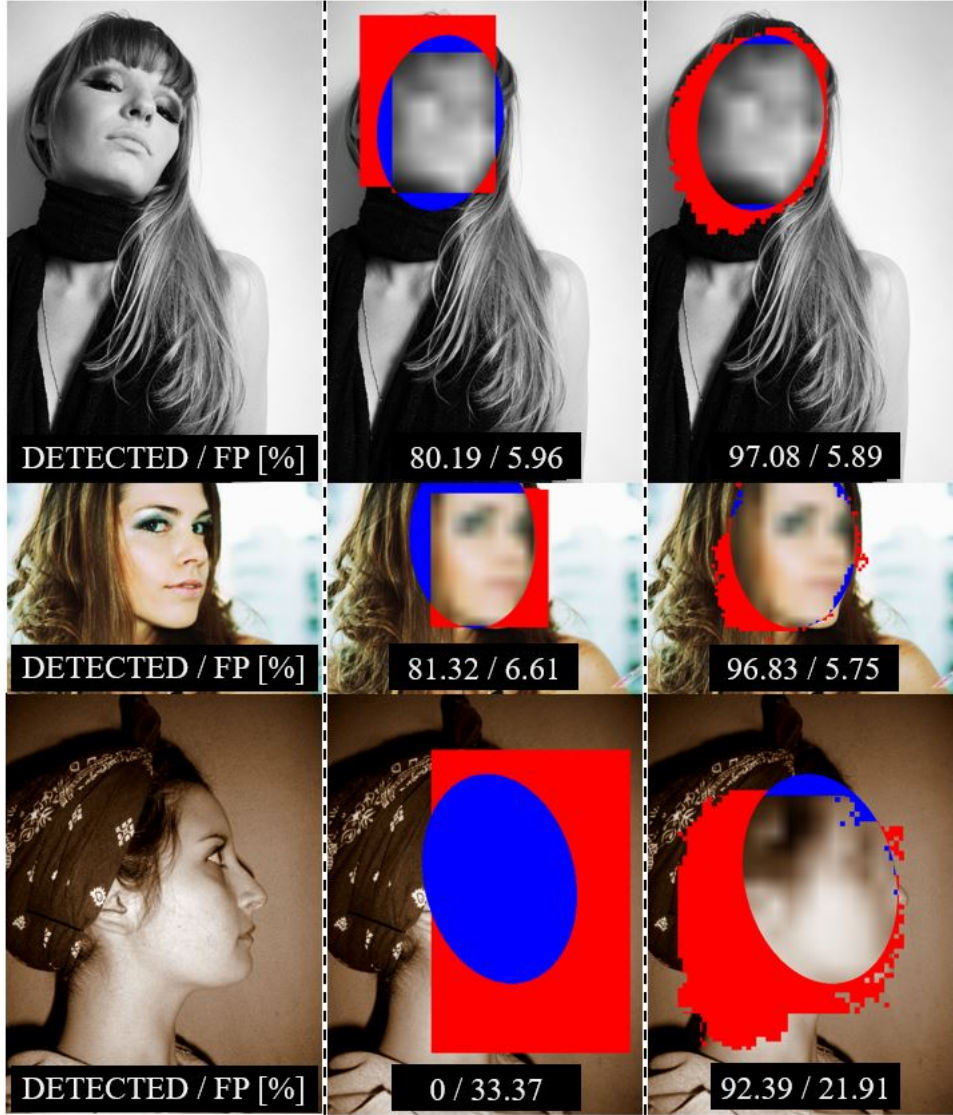


Figure 19: Results comparing our proposed method with the results of the approach presented in Hu et al. for three different example images.

Therefore, as seen from the figure, the proposed algorithm preserves the face shapes better than the approach presented in [42]. Moreover, as shown in the bottom image of Fig. 19, the proposed algorithm is particularly superior to the existing approach when faces are largely rotated in the yaw direction.

## B. Robot Control Evaluation

### • Robot Hardware Specification

We performed our experiments on a TurtleBot3 Burger which is a small, affordable, and programmable ROS-based mobile robot. The setup of the robot used in our experiments is shown in Fig. 20. In order to obtain a wider viewing angle, a stereo ZED camera was installed at 1.32m height. A Raspberry Pi 3 as a motor control computer and an Nvidia

Xavier as a companion computer were mounted on the robot.

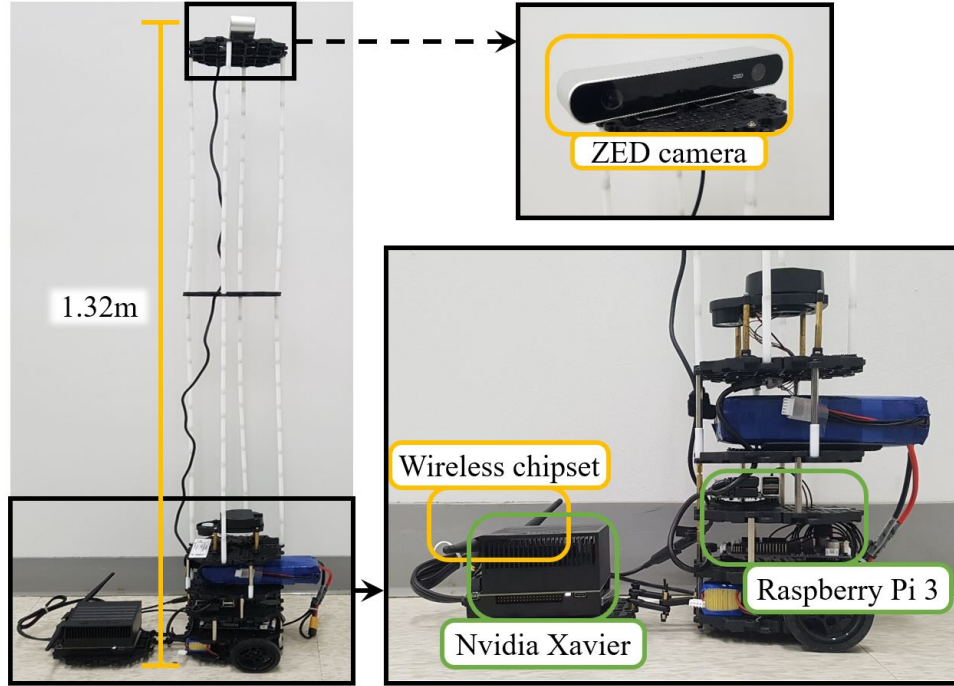


Figure 20: Face detection robot used in our experiments

- **Experiments of Face Detection Network with Nvidia Xavier**

Because the Nvidia Xavier module is equipped with a 512-core Volta GPU with Tensor cores, it is suitable for large-scale matrix operations which are needed for efficient neural networks computation. However, since the combined CPU and GPU memory is limited by only 16 GB and because its power consumption is relatively high, it is difficult to operate our face detection algorithm by simply porting the program to the Nvidia Xavier on the robot. To overcome these problems, an SSD memory card with a capacity of 256GB was additionally installed for a swap memory, and a 4-cell LiPo battery with a voltage of 14.8V was installed for power supply. In order to minimize the amount of computations, IOU threshold of NMS (non maximum suppression) is set to 0.1 so that not many bounding boxes are drawn on faces. Since the run-time is 2.2FPS when running in an Nvidia Xavier, we can utilize our face detection algorithm to any other real-time application such as SLAM for the TurtleBot3 Burger with a maximum speed of 0.26m/s.

- **ORB-SLAM2**

The first step for ORB SLAM2 is the extraction of several features from an image such as edges, corners and lines [47]. Fig. 21 shows the comparison of the feature extraction results at various resolutions. The features are well extracted even for small eyes and noses in HR images of  $720 \times 1280$  pixels. However, the extraction algorithm fails to extract any feature in extreme LR images of  $24 \times 32$  pixels as shown in Fig. 21 (b).

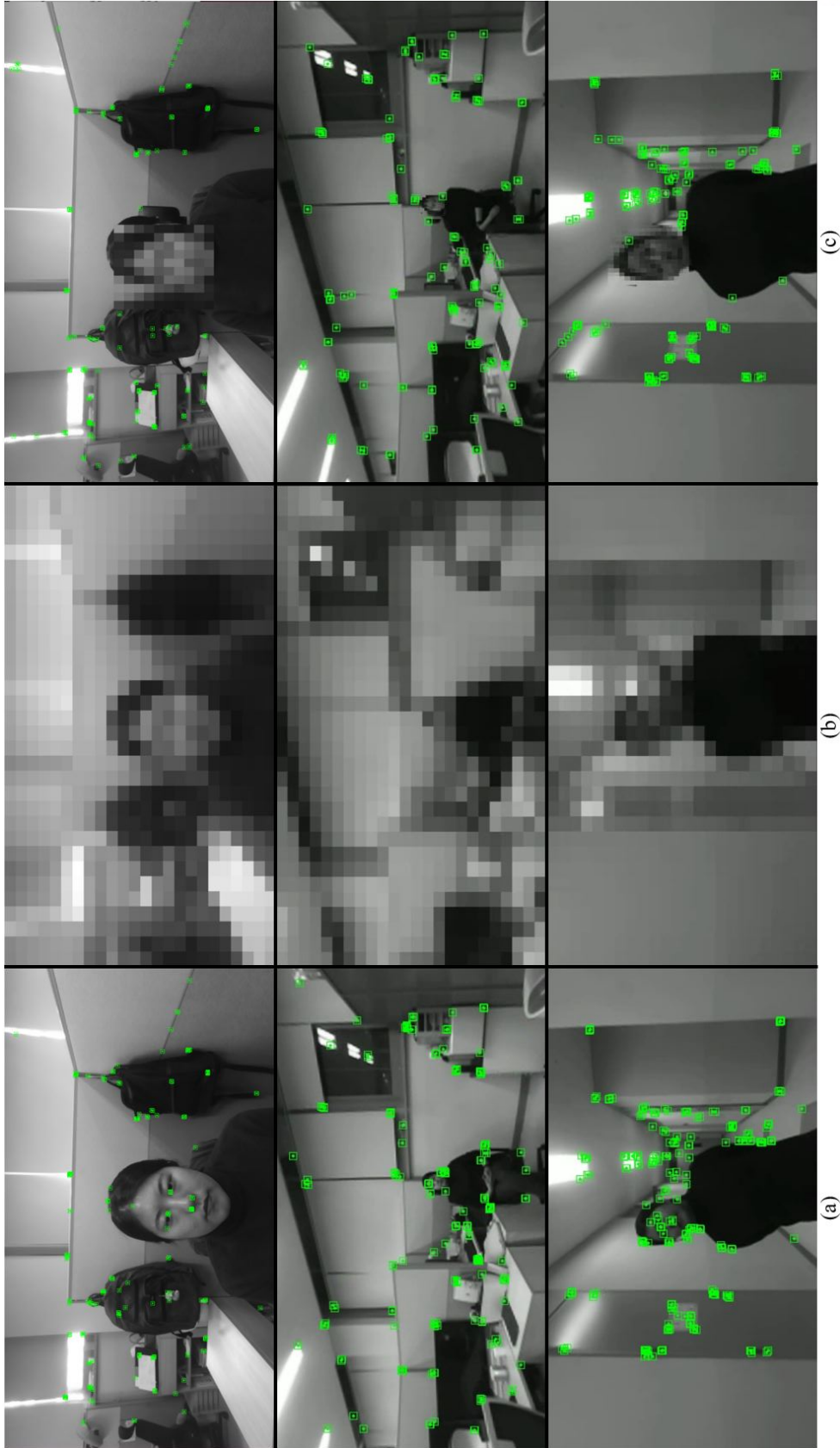


Figure 21: Comparison of the feature extraction results at various resolutions.



The feature points extracted from our privacy preserved image in Fig. 21 (c) are almost the same as the feature points in Fig. 21 (a) except for a few features in the face part. Therefore, with our proposed system, a robot can still operate SLAM smoothly even without any aid from or connection to the ground station. Video results for feature extraction and ORB-SLAM2 can be found in the following link: [https://youtu.be/\\_W6e6xPRsM0](https://youtu.be/_W6e6xPRsM0).

## Conclusion

We developed a patrol robot system proposing a novel privacy preserving face detection camera system. Using the proposed system, every snapshot is transformed to an image with faces of extreme LR on a real-time basis without any computational aid from the ground station. Since any of the faces is always in extreme LR even during the face detection process, privacy of the people in the snapshot can be fundamentally preserved. We experimentally confirmed that a robot with the proposed privacy-preserved robot vision can perform its real-time operation such as SLAM. Upon request from users, the face detection for privacy protection can be extended to a face and body detection.

## References

- [1] M. U. Kim and H. J. Yang, "Min-SINR maximization with DL SWIPT and UL WPCN in multi-antenna interference networks," *IEEE Wireless Communications Letters*, vol. 6, no. 3, pp. 318–321, 2017.
- [2] M. U. Kim and H. Jong Yang, "RNN-based node selection for sensor networks with energy harvesting," in *2018 International Conference on Information and Communication Technology Convergence (ICTC)*, pp. 1316–1318, 2018.
- [3] M. U. Kim, H. Lee, H. J. Yang, and M. S. Ryoo, "Privacy-preserving robot vision with anonymized faces by extreme low resolution," in *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 462–467, 2019.
- [4] "Scenarios and requirements for small cell enhancements," 3GPP TR 36.932, v12.1.0.
- [5] "Small cell enhancements for E-UTRA and E-UTRAN - physical layer aspects," 3GPP TR 36.872, v12.1.0.
- [6] "Increasing capacity in wireless broadcast systems using distributed transmission/directional reception (DTDR)," US Patent 5345599.
- [7] C. K. Au-yeung and D. J. Love, "On the performance of random vector quantization limited feedback beamforming in a MISO system," *IEEE Trans. on Wireless Commun.*, vol. 6, no. 2, pp. 458–462, 2007.
- [8] H. Wang, X. Zhou, and M. C. Reed, "Analytical evaluation of coverage-oriented femtocell network deployment," in *2013 IEEE Int'l Conf. on Commun. (ICC)*, pp. 5974–5979, 2013.
- [9] H. J. Yang, W. Shin, B. C. Jung, and A. Paulraj, "Opportunistic interference alignment for MIMO interfering multiple-access channels," *IEEE Trans. on Wireless Commun.*, vol. 12, no. 5, pp. 2180–2192, 2013.
- [10] H. J. Yang, W. Shin, B. C. Jung, C. Suh, and A. Paulraj, "Opportunistic downlink interference alignment for multi-cell MIMO networks," *IEEE Trans. on Wireless Commun.*, vol. 16, no. 3, pp. 1533–1548, 2017.
- [11] R. V. Nee and R. Prasad, *OFDM for Wireless Multimedia Communications*. Artech House, Inc., January 2000.

- [12] D. Tse and P. Viswanath, *Fundamentals of Wireless Communication*. Cambridge University Press, 2005.
- [13] D. Lopez-Perez, I. Guvenc, G. de la Roche, M. Kountouris, T. Q. S. Quek, and J. Zhang, "Enhanced intercell interference coordination challenges in heterogeneous networks," *IEEE Wireless Commun.*, vol. 18, no. 3, pp. 22–30, 2011.
- [14] D. Gesbert, S. Hanly, H. Huang, S. Shamai Shitz, O. Simeone, and W. Yu, "Multi-cell MIMO cooperative networks: A new look at interference," *IEEE J. Selected Areas in Commun.*, vol. 28, no. 9, pp. 1380–1408, 2010.
- [15] M. Rahimi, H. Shah, G. S. Sukhatme, J. Heideman, and D. Estrin, "Studying the feasibility of energy harvesting in a mobile sensor network," in *IEEE Int'l Conf. on Robotics and Automntion*, vol. 1, pp. 19–24, September 2003.
- [16] S. Priya and D. J. Inman, *Energy Harvesting Technologies*. Springer, 2009.
- [17] H. S. Dhillon, Y. Li, P. Nuggehalli, Z. Pi, and J. G. Andrews, "Fundamentals of heterogeneous cellular networks with energy harvesting," *IEEE Trans. Wireless Commun.*, vol. 13, pp. 2782–2797, May 2014.
- [18] R. Zhang and C. K. Ho, "MIMO broadcasting for simultaneous wireless information and power transfer," *IEEE Trans. Wireless Commun.*, vol. 12, pp. 1989–2001, May 2013.
- [19] L. Zhao, X. Wang, and K. Zheng, "Downlink hybrid information and energy transfer with massive MIMO," *IEEE Trans. Wireless Commun.*, vol. 15, pp. 1309–1322, February 2016.
- [20] J. Xu, L. Liu, and R. Zhang, "Multiuser MISO beamforming for simultaneous wireless information and power transfer," *IEEE Trans. Signal Process.*, vol. 62, pp. 4798–4810, September 2014.
- [21] D. W. K. Ng, E. S. Lo, and R. Schober, "Robust beamforming for secure communication in systems with wireless information and power transfer," *IEEE Trans. Wireless Commun.*, vol. 13, pp. 4599–4615, August 2014.
- [22] A. Ghazanfari, H. Tabassum, and E. Hossain, "Ambient RF energy harvesting in ultra-dense small cell networks : Performance and trade-offs," *IEEE Wireless Commun.*, vol. 23, pp. 38–45, 2016.
- [23] H. Ju and R. Zhang, "Throughput maximization in wireless powered communication networks," *IEEE Trans. Wireless Commun.*, vol. 13, pp. 418–428, January 2014.
- [24] G. Yang, C. K. Ho, R. Zhang, and Y. Liang Guan, "Throughput optimization for massive MIMO systems powered by wireless energy transfer," *IEEE J. Selected Areas in Commun.*, vol. 33, pp. 1640–1650, August 2015.

- [25] H. Lee, K.-J. Lee, H.-B. Kong, and I. Lee, "Sum rate maximization for multi-user MIMO wireless powered communication networks," *IEEE Trans. Vehicular Technology*, 2016.
- [26] T. Le, K. Mayaram, and T. Fiez, "Efficient far-field radio frequency energy harvesting for passively powered sensor networks," *IEEE Journal of Solid-State Circuits*, vol. 43, no. 5, pp. 1287–1302, 2008.
- [27] M. Schubert and H. Boche, "Iterative multiuser uplink and downlink beamforming under SINR constraints," *IEEE Trans. Signal Processing*, vol. 53, pp. 2324–2334, July 2005.
- [28] L. Vandenberghe and S. Boyd, "Semidefinite programming," *Society for Industrial and Applied Mathematics (SIAM) Review*, vol. 38, pp. 49–95, March 1996.
- [29] M. C. Grant, S. Boyd, and Y. Ye, "CVX: Matlab software for disciplined convex programming," June 2015.
- [30] Z.-Q. Luo, W.-K. Ma, A. M.-C. So, Y. Ye, and S. Zhang, "Semidefinite relaxation of quadratic optimization problems," *IEEE Signal Processing Magazine*, vol. 27, pp. 20–34, May 2010.
- [31] J. Yang, X. Wu, and J. Wu, "Optimal scheduling of collaborative sensing in energy harvesting sensor networks," *IEEE J. Sel. Areas Commun.*, pp. 512–523, March 2015.
- [32] F. Shan, J. Luo, W. Wu, and X. Shen, "Optimal wireless power transfer scheduling for delay minimization," in *IEEE INFOCOM*, April 2016.
- [33] M. U. Kim and H. J. Yang, "Min-sinr maximization with dl swipt and ul wpcn in multicell multi-antenna networks," *IEEE Wireless Commun. Lett.*, vol. 6, pp. 318–321, June 2017.
- [34] "Ieee standard for low-rate wireless networks," *IEEE Std 802.15.4-2015 (Revision of IEEE Std 802.15.4-2011)*, April 2016.
- [35] "Ieee standard for local and metropolitan area networks—part 15.4: Low-rate wireless personal area networks (lr-wpans) amendment 3: Physical layer (phy) specifications for low-data-rate, wireless, smart metering utility networks," *IEEE Standard 802.15.4g-2012*, April 2012.
- [36] E. Zeng, S. Mare and F. Roesner, "End user security and privacy concerns with smart homes," in *Proc. 13th Symp. on Usable Privacy and Security*, pp. 65–80, 2017.
- [37] R. Templeman, M. Korayem, D. Crandall, and A. Kapadia, "Placeavoider: Steering first-person cameras away from sensitive spaces," in *Proc. Network and Distributed System Security Symp.*, 2014.
- [38] Z. Ren, Y. J. Lee and M. S. Ryoo, "Learning to anonymize faces for privacy preserving action detection," in *Proc. ECCV*, pp. 620–636, 2018.

- [39] M. S. Ryoo, B. Rothrock, C. Fleming, and H. J. Yang, “Privacy-preserving human activity recognition from extreme low resolution,” in *Proc. the 31st AAAI*, pp. 4255–4262, 2017.
- [40] M. S. Ryoo, K. Kim, and H. J. Yang, “Extreme low resolution activity recognition with multi-siamese embedding learning,” in *Proc. the 32nd AAAI*, 2018.
- [41] Z. Wu, Z. Wang, Z. Wang and H. Jin, “Towards privacy-preserving visual recognition via adversarial training: A pilot study,” in *Proc. ECCV*, pp. 627–645, 2018.
- [42] P. Hu, and D. Ramanan, “Finding tiny faces,” in *Proc. CVPR*, pp. 1522–1530, 2017.
- [43] Y. Bai, Y. Zhang, M. Ding and B. Ghanem, “Finding tiny faces in the wild with generative adversarial network,” in *Proc. CVPR*, pp. 21–30, 2018.
- [44] D. H. Kim, W. Yun and J. Lee, “Tiny frontal face detection for robots,” in *Proc. Int’l Conf. on Human-Centric Computing*, pp. 1–4, 2010.
- [45] A. A. Efros, A. C. Berg, G. Mori and J. Malik, “Recognizing action at a distance,” in *Proc. ICCV*, pp. 726—733, 2003.
- [46] S. Ali, A. Basharat and M. Shah, “Chaotic invariants for human action recognition,” in *Proc. ICCV*, pp. 1–8, 2007.
- [47] R. Mur-Artal and J. D. Tardós, “ORB-SLAM2: an open-source SLAM system for monocular, stereo and RGB-D cameras,” *IEEE Trans. Robot*, vol. 33, no. 5, pp. 1255–1262, 2017.
- [48] M. Koestinger, P. Wohlhart, P. M. Roth, and H. Bischof, “Annotated facial landmarks in the wild: A large-scale, real-world database for facial landmark localization,” in *Proc. First IEEE Int’l Workshop on Benchmarking Facial Image Analysis Technologies*, 2011.
- [49] E. Shelhamer, J. Long, and T. Darrell, “Fully convolutional networks for semantic segmentation,” *IEEE Trans. Pattern Anal. Match. Intell.*, vol. 39, no. 4, pp. 640–651, 2017.
- [50] K. Simonyan, and A. Zisserman, “Very deep convolutional networks for large-scale image recognition,” in *Proc. IEEE Int’l Conf. on Learning Representations*, 2015.
- [51] S. Ren, K. He, R. B. Girshick and J. Sun, “Faster R-CNN: towards real-time object detection with region proposal networks,” *CoRR*, vol. abs/1506.01497, 2015.
- [52] V. Jain, and E. Learned-Miller, “Fdldb: A benchmark for face detection in unconstrained settings,” Tech. Rep. UM-CS-2010-009, University of Massachusetts, Amherst, 2010.

## Acknowledgements

I am deeply grateful to those who helped me a lot in my doctoral studies. First, my advisor, Professor Hyun Jong Yang, is my biggest supporter. He has advised me passionately to grow as a good researcher and allowed me to study in a good environment. Thanks to his supervision, despite several obstacles, I was able to immerse myself in research enthusiastically and produce good results. In addition, as a senior in life, he allowed me to indirectly experience various things I had never experienced. Once again, I would like to express my sincere gratitude to Professor Hyun Jong Yang. I am also grateful to the members of the AIS lab for sharing good studies with me. Thanks to working with them, I was able to learn what co-work is and how to get better results. As well, I would like to thank my friends, Sunjung Kang and Hye-young Shin, and UNIST Rowing Club members, especially Sae-eun Kim, and Tae-joon Kim, those who have always helped me during my graduate school life.

Finally, I would like to express my deep gratitude to my family, especially my parents. Without their strong support, I would not have been able to study for a long time.

In the future, I will try to live a life that rewards the many loves I have received so far. Again, thanks to everyone.

